



US005469354A

**United States Patent** [19]

Hatakeyama et al.

[11] **Patent Number:** 5,469,354[45] **Date of Patent:** Nov. 21, 1995[54] **DOCUMENT DATA PROCESSING METHOD  
AND APPARATUS FOR DOCUMENT  
RETRIEVAL**5-55912 8/1993 Japan.  
5-76068 10/1993 Japan.  
WO90/16036 2/1990 WIPO.[75] **Inventors:** Atsushi Hatakeyama, Kokubunji;  
Hiromichi Fujisawa; Kanji Kato, both  
of Tokorozawa; Hisamitsu Kawaguchi,  
Sagamihara; Naoki Minegishi, Osaka;  
Katsumi Tada, Kokubunji; Satoshi  
Asakawa, Hirakata, all of Japan[73] **Assignee:** Hitachi, Ltd., Tokyo, Japan[21] **Appl. No.:** 843,162[22] **Filed:** Feb. 28, 1992**Related U.S. Application Data**[63] Continuation-in-part of Ser. No. 555,483, Aug. 9, 1990, Pat.  
No. 5,168,533.[30] **Foreign Application Priority Data**Feb. 28, 1991 [JP] Japan ..... 3-058311  
Dec. 25, 1991 [JP] Japan ..... 3-342695[51] **Int. Cl.<sup>6</sup>** ..... G06F 17/21[52] **U.S. Cl.** ..... 364/419.19; 364/419.13;  
364/419.07; 364/225.3[58] **Field of Search** ..... 364/419.19, 419.13,  
364/419.07, 225.3; 395/600[56] **References Cited****U.S. PATENT DOCUMENTS**4,870,568 9/1989 Kahle et al. .... 364/225.3  
5,051,947 9/1991 Messenger et al. .... 364/956.1  
5,168,533 12/1992 Kato et al. .... 364/200  
5,206,949 4/1993 Cochran et al. .... 395/600  
5,220,625 6/1993 Hatakeyama et al. .... 382/54**FOREIGN PATENT DOCUMENTS**0437615A1 7/1991 European Pat. Off. .  
63-198124 8/1988 Japan .  
3-125263 5/1991 Japan .**OTHER PUBLICATIONS**Mukhopadhyay et al., An Intelligent System for Document  
Retrieval in Distributed Office Environments, Journal of the  
American Society for Information Science, Jun. 17, 1985.  
"State Machines Find the Pattern", System Design/Soft-  
ware, 8167 Computer Design, May 1985, No. 5, Littleton,  
Mass.**Primary Examiner**—Robert A. Weinhardt**Assistant Examiner**—Frantzy Poinvil**Attorney, Agent, or Firm**—Antonelli, Terry, Stout & Kraus[57] **ABSTRACT**

High-speed full document retrieval method and system capable of providing result of retrieval within practically acceptable short search time. Upon registration of documents in a document database, condensed texts are created by decomposing each of textual character strings of the documents to be registered into fragmental character strings in dependence on character species and by checking mutual inclusion relations existing among the fragmental character strings. A component character table is created in which characters occurring in each of the condensed texts are registered without duplication. The condensed texts and the component character table are registered in the data base together with the texts of the documents to be registered. Upon retrieval of a document containing a search term designated by a user, a component character table search is first executed to extract those documents which contain all species of characters constituting the search term by consulting the component character table, and subsequently a condensed text search is executed by consulting the condensed texts of the documents. Finally, a text body search is executed for extracting a document which satisfies query condition imposed on the search term by consulting the texts of the documents extracted through the component character table search and the condensed text search.

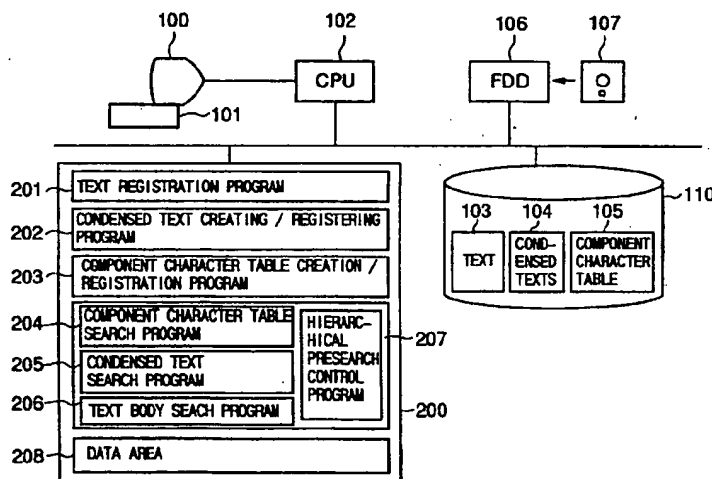
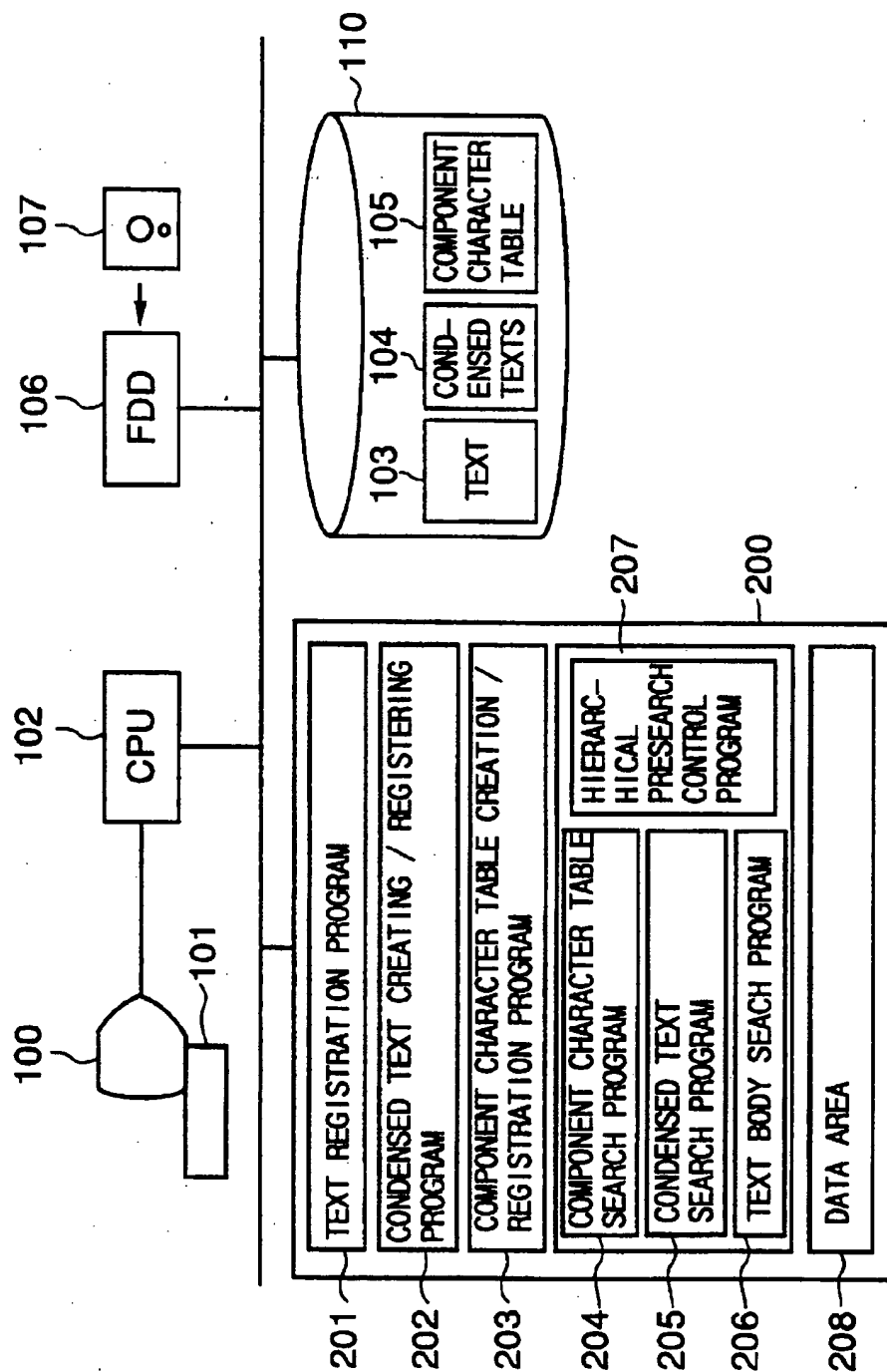
**50 Claims, 66 Drawing Sheets**

FIG. 1



## FIG. 2

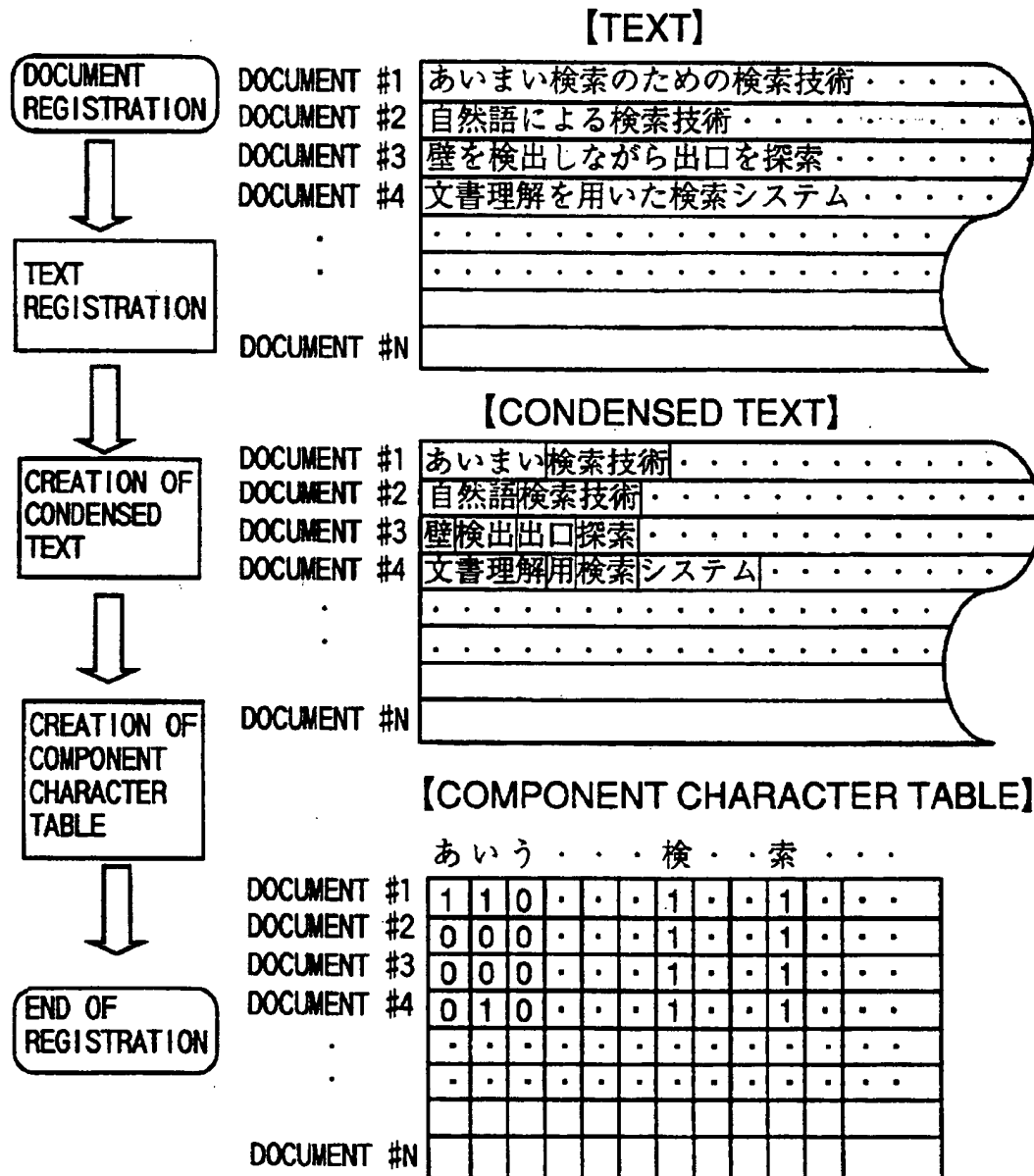


FIG. 3

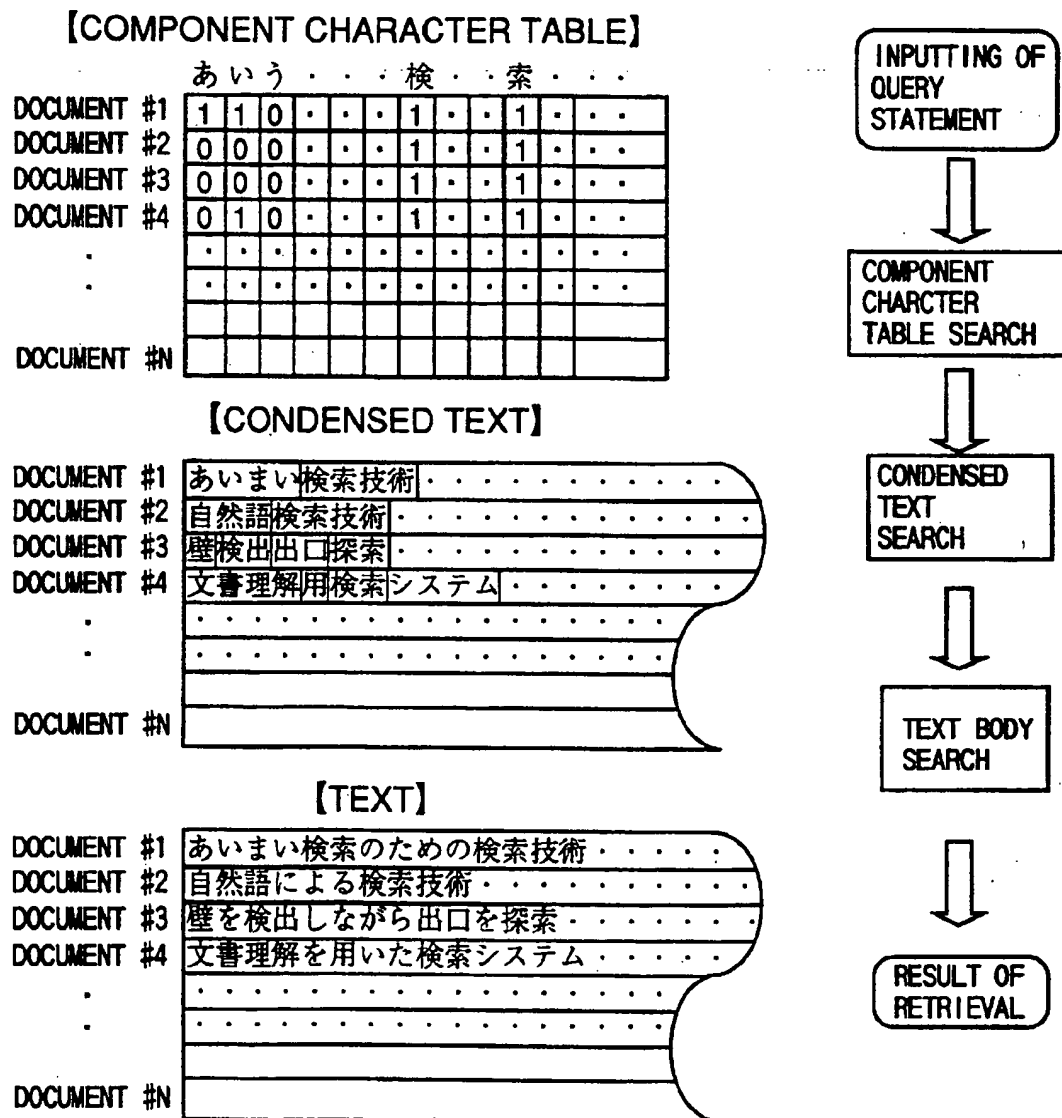
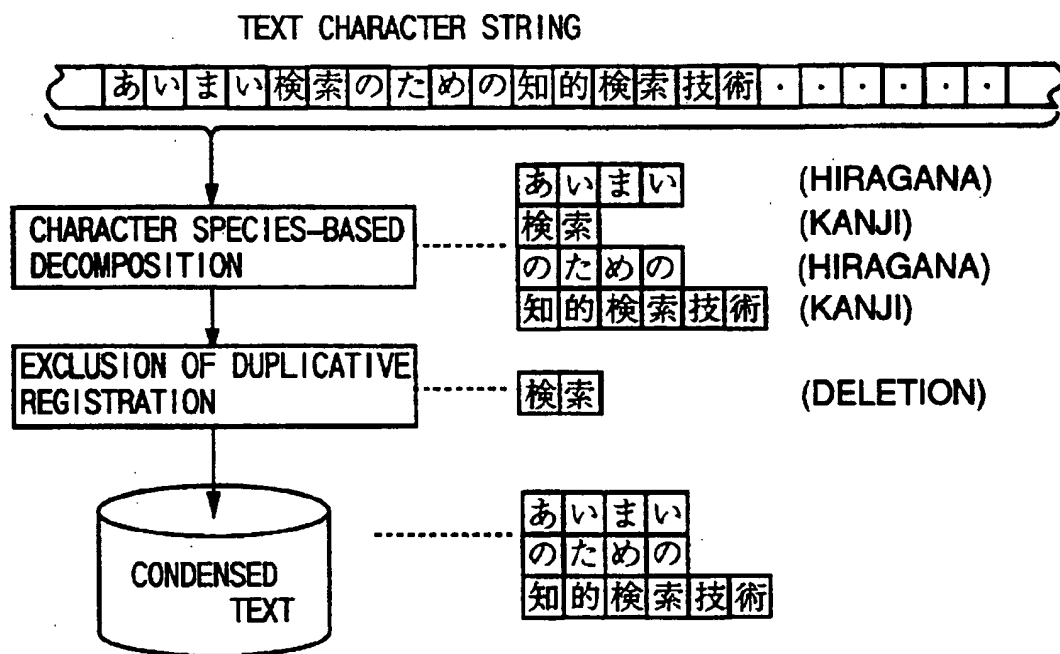


FIG. 4



## FIG. 5

DOCUMENT #1	あいまい,のための,検索技術
DOCUMENT #2	自然語,による,検索技術
DOCUMENT #3	壁,を,検出,しながら,出口,検索
DOCUMENT #4	文書理解,を,用,いた,検索,システム
.	
.	
DOCUMENT #N	

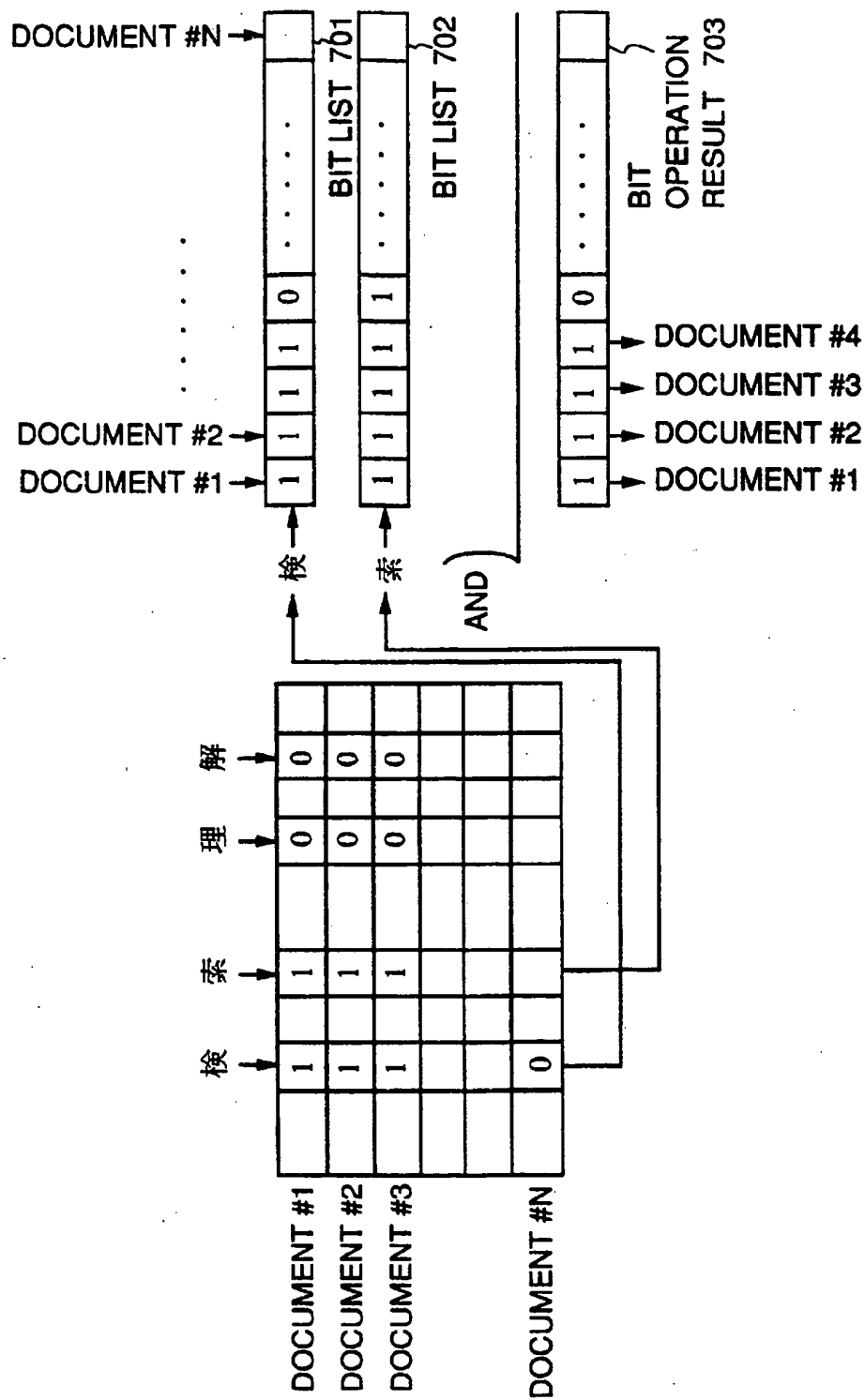
FIG. 6

CHARACTER CODE	(0000)H	(SPACE) (8140)H	あ (82A0)H	い (82A2)H	ま (82DC)H	を (82F0)H	垂 (82F0)H	技 (889F)H	検 (8B5A)H	索 (8C9F)H	論 (EA9E)H
DOCUMENT #1	00	...	1	...	1	...	0	...	1	...	0
DOCUMENT #2	00	...	1	...	1	...	0	...	1	...	0
DOCUMENT #3	00	...	1	...	1	...	0	...	1	...	0
DOCUMENT #4	00	...	1	...	1	...	0	...	1	...	0
.	00	...	1	...	1	...	0	...	1	...	0
.	00	...	1	...	1	...	0	...	1	...	0
DOCUMENT #N											

ENTRY ID  
NUMBER

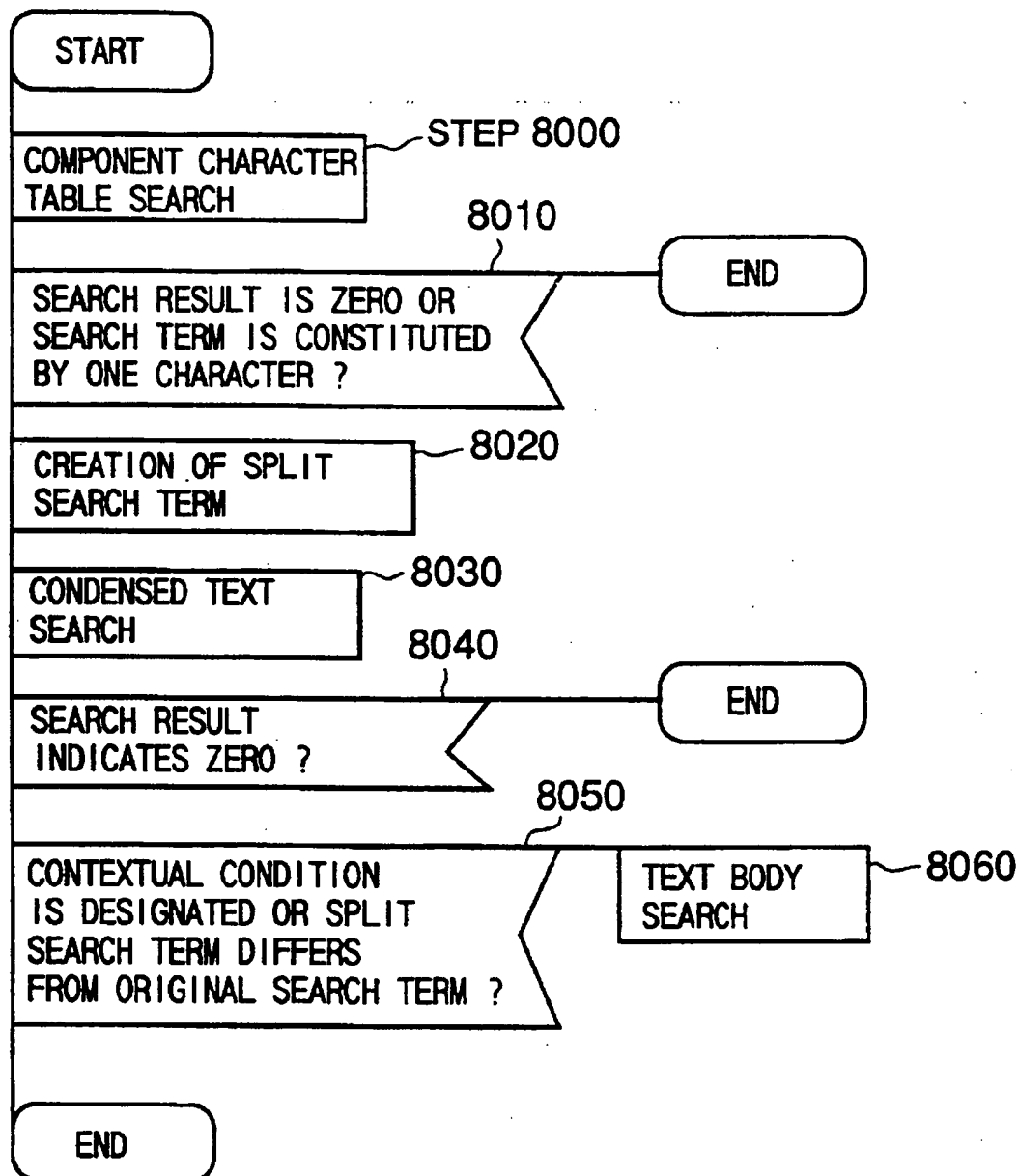
01 .... 33088 .... 33440 33442 .. 33500 .. 33520 .. 34975 .. 35674 .. 35999 .. 36341 .. 60062

FIG. 7

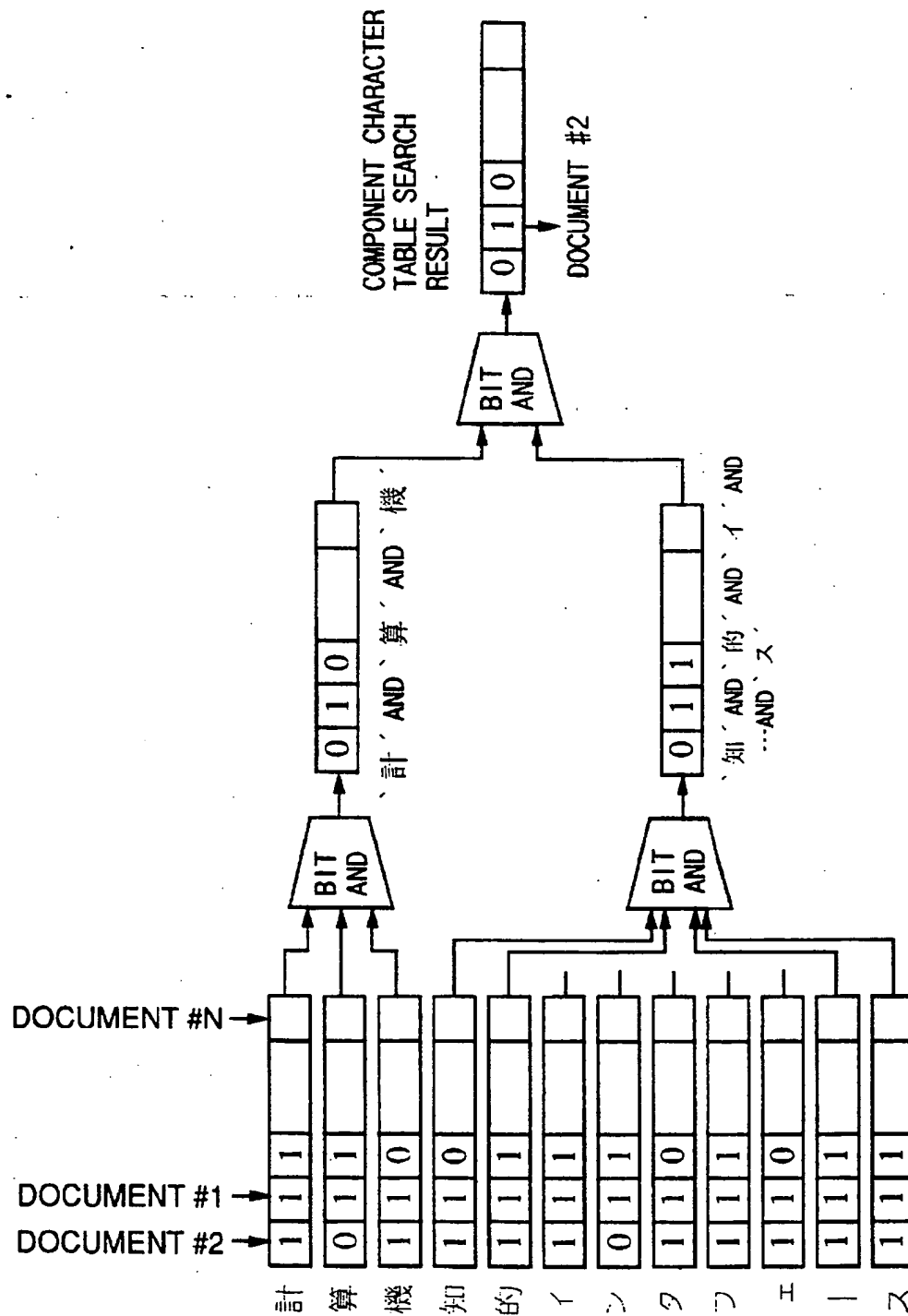




## FIG. 8



9. 611



## FIG. 10

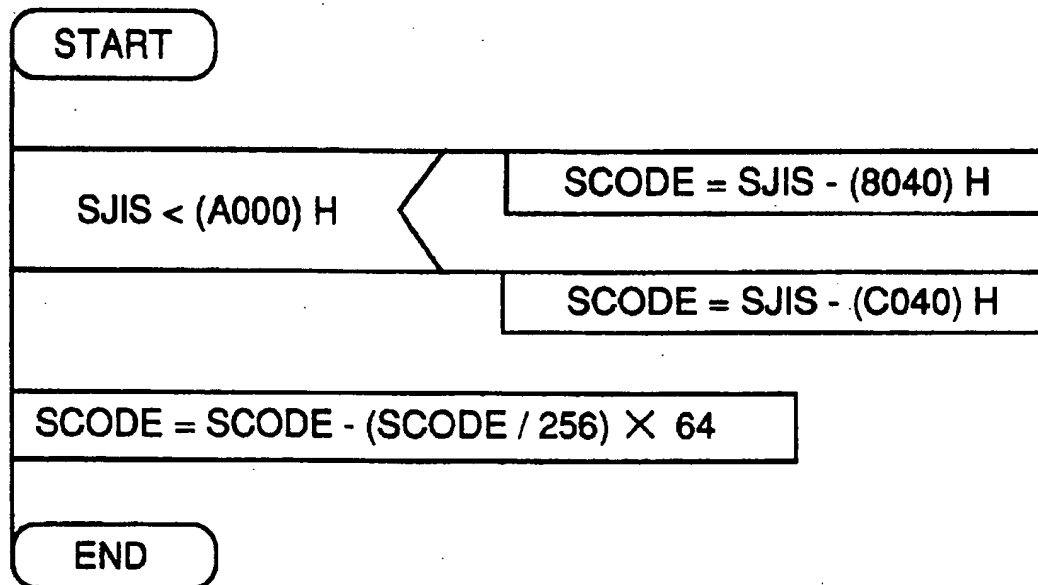


FIG. 11

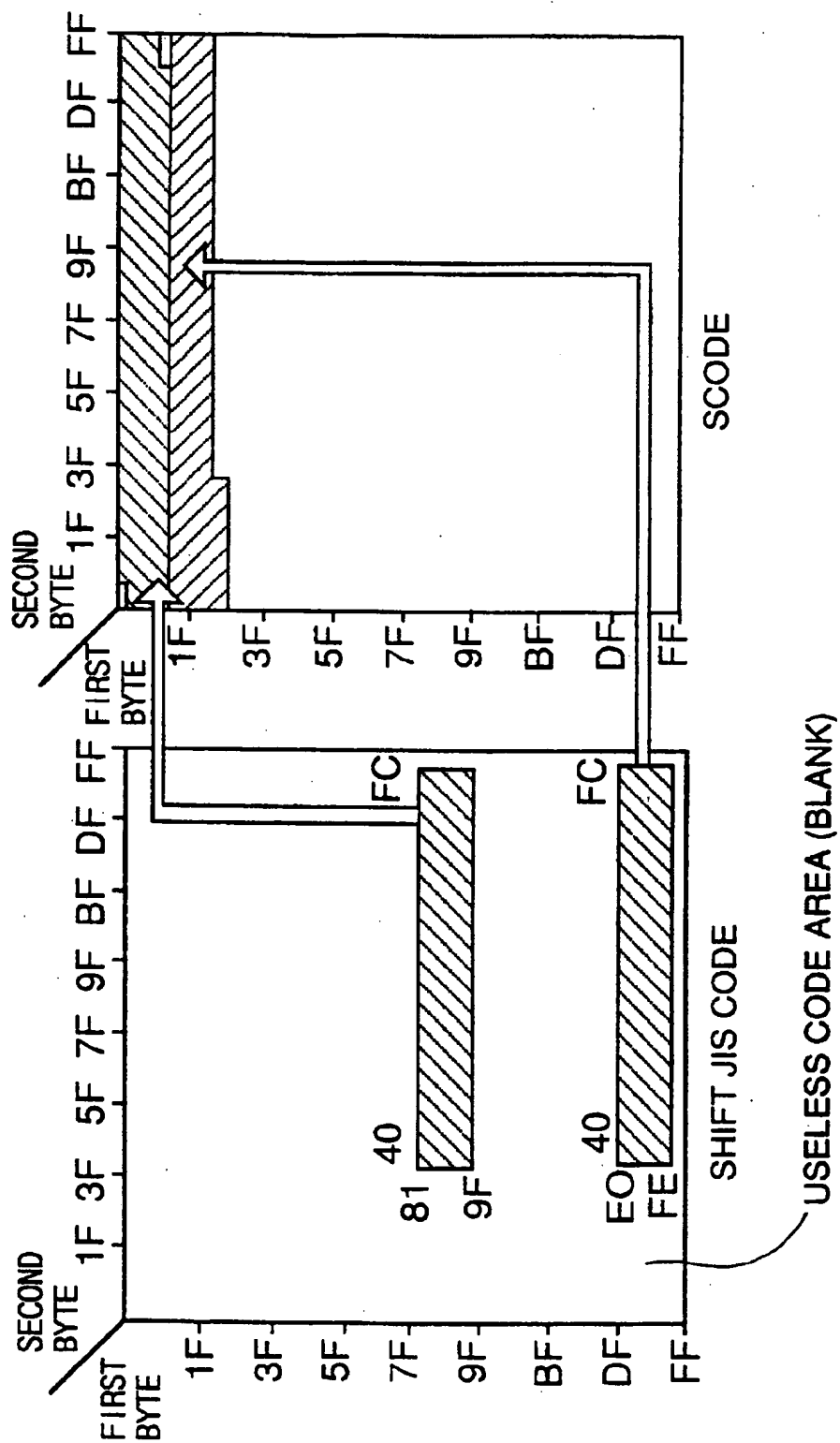


FIG. 12

TRANSFORMED CHARACTER CODE	(0000)H	(SPACE) (00C0)H	あ (01E0)H	い (01E2)H	ま (021C)H	を (0230)H	重 (065F)H	技 (085A)H	検 (095F)H	索 (0A75)H	備 (1FDE)H
DOCUMENT #1	00	...	1	01	...	0	...	1	...	1	...
DOCUMENT #2	00	...	0	00	...	0	...	1	...	1	...
DOCUMENT #3	00	...	0	00	...	1	...	0	...	1	...
DOCUMENT #4	00	...	0	01	...	1	...	0	...	1	...
.	00	...	1	...	...	...	...	...	...	...	...
DOCUMENT #N	...	...	...	...	...	...	...	...	...	...	...

ENTRY ID    0 1    ...    192    ...    480 482    ..    540    ..    560    ..    1631    ..    2138    ..    2399    ..    2677    ..    8158  
NUMBER

FIG. 13

HASHED CHARACTER CODE	ま (0)	を (28)	を (48)	技 (90)	理 (95)	索 (117)	索 (SPACE) (182)	検 (351)	あ (480)	い (482)	
DOCUMENT #1	0 0	1 ...	0 ...	1 ...	0 ...	1 ...	1 ...	1 ...	1 ...	1 0 1	...
DOCUMENT #2	0 0	0 ...	0 ...	1 ...	0 ...	1 ...	1 ...	1 ...	0 ...	0 0 0	...
DOCUMENT #3	0 0	0 ...	1 ...	0 ...	0 ...	1 ...	1 ...	1 ...	0 ...	0 0 0	...
DOCUMENT #4	0 0	0 ...	1 ...	0 ...	0 ...	1 ...	1 ...	1 ...	0 ...	0 0 1	...
.	0 0	...									
.											
DOCUMENT #N											

ENTRY ID    0 1 ... 28 ... 48 ... 90 .. 95 ... 117 ... 192 ... 351 ... 480 482 ... 511

NUMBER

## FIG. 14

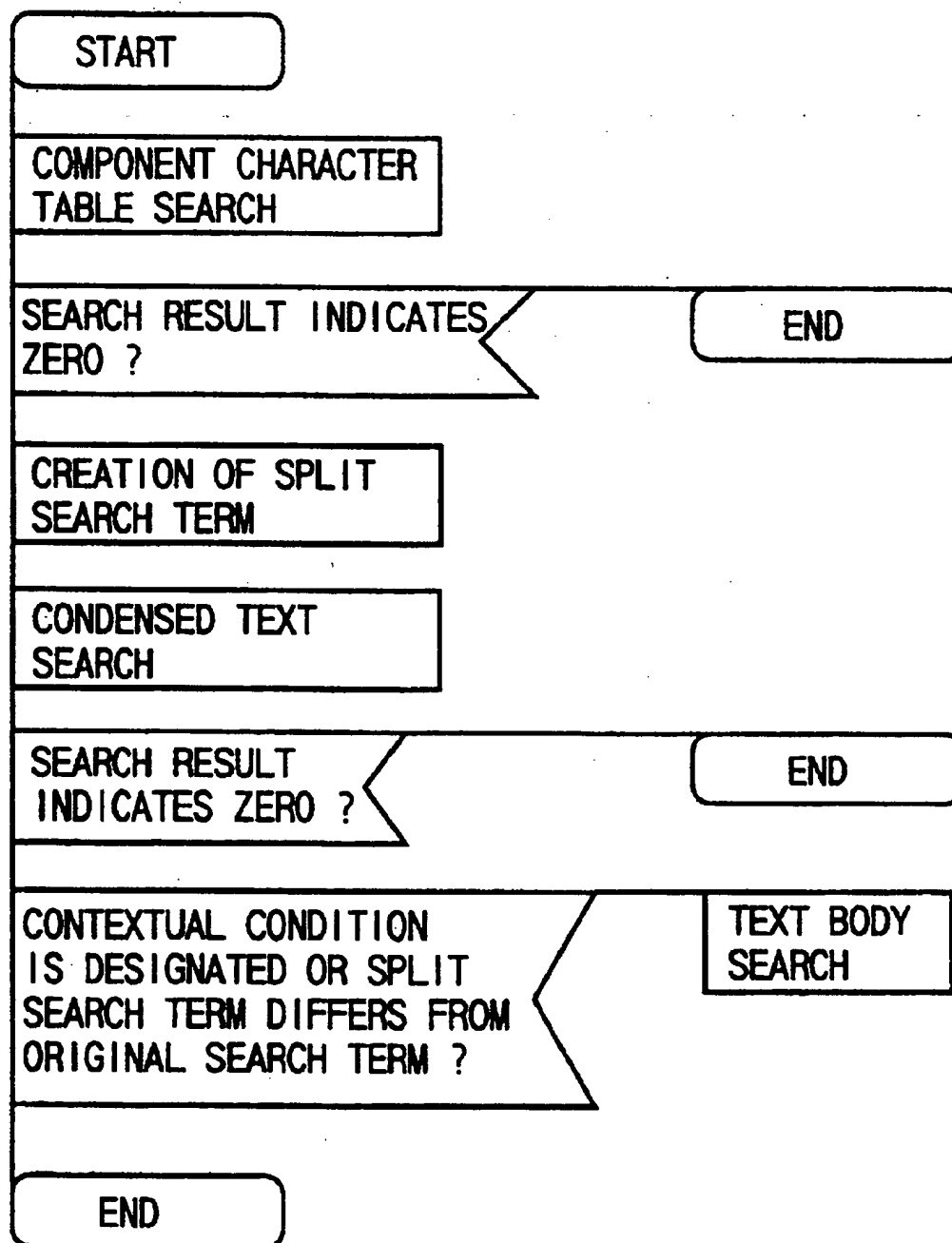
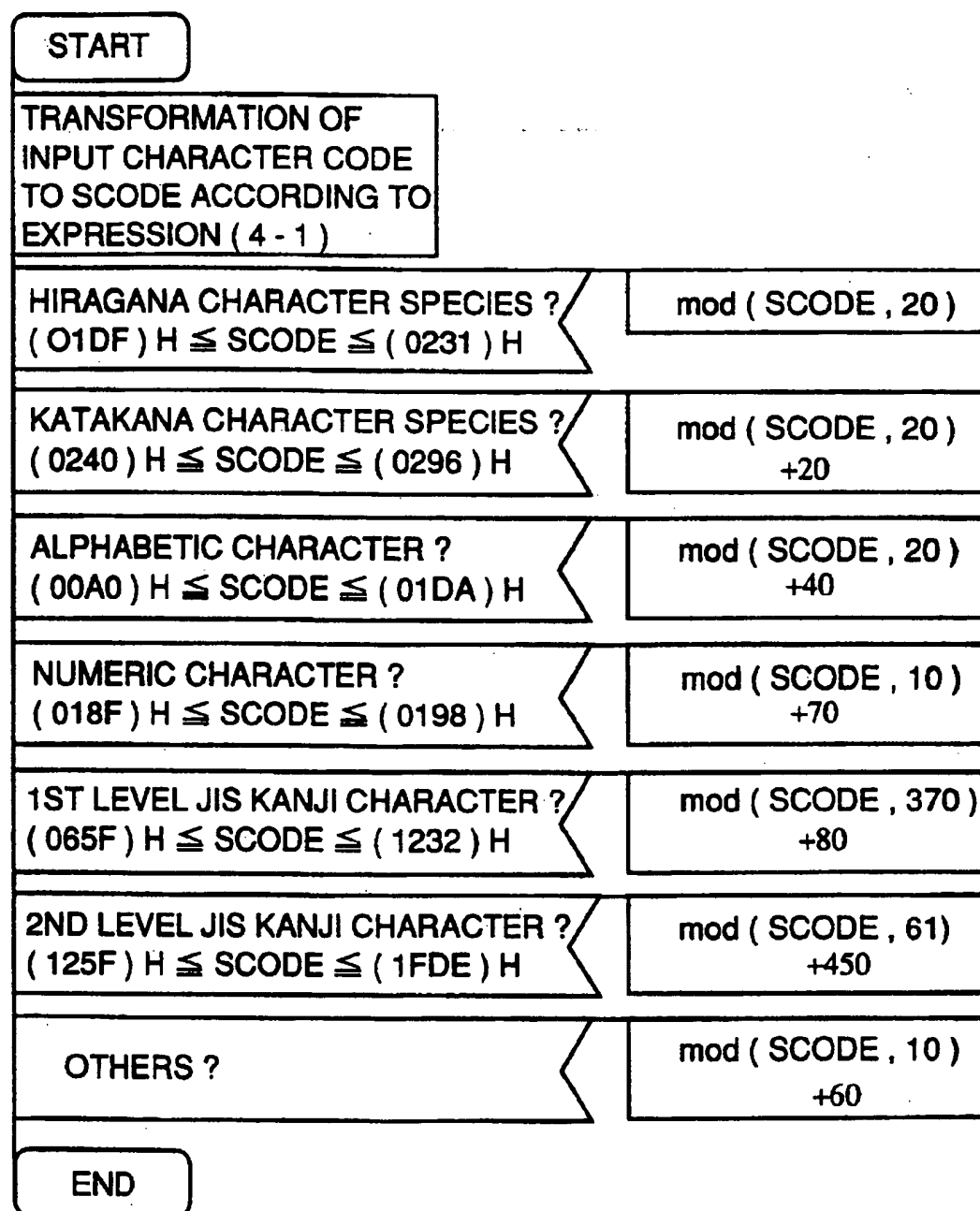


FIG. 15

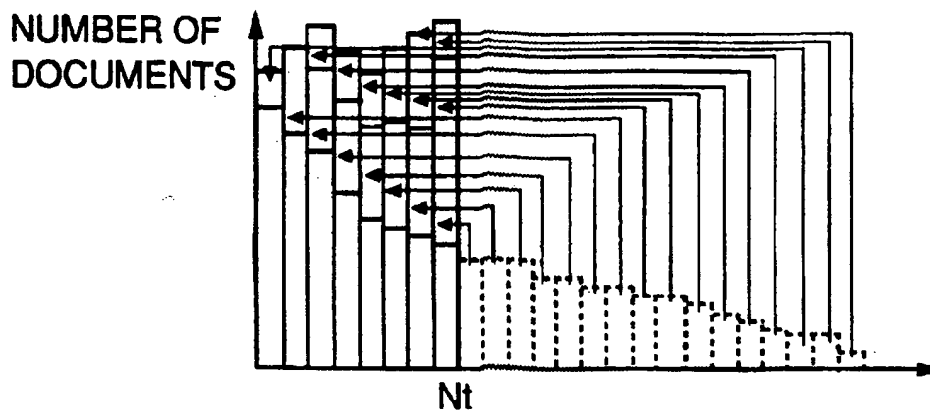
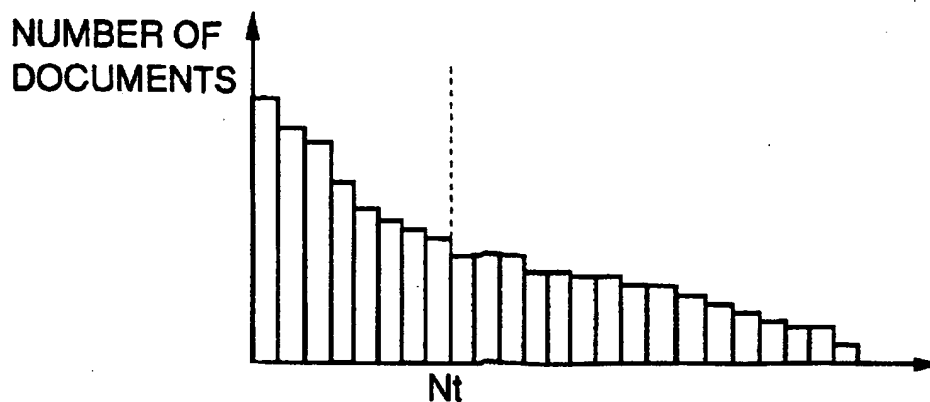
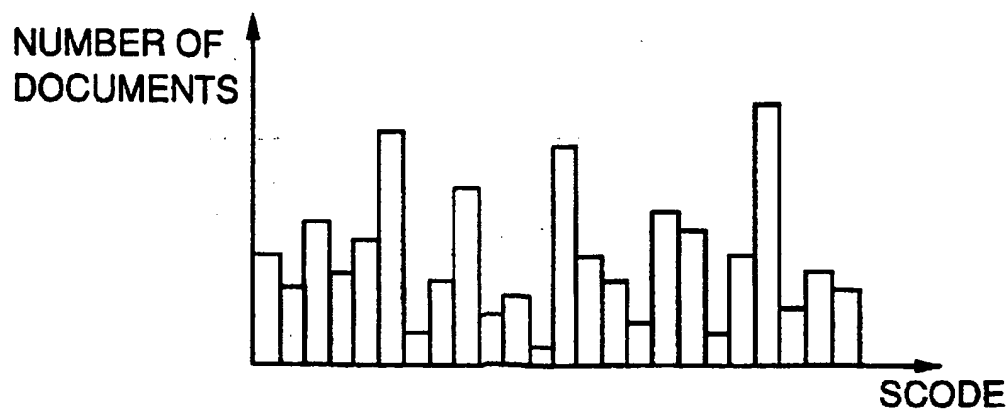
	あ	い	ま	を	(SPACE)	技	検	索																
DOCUMENT #1	1	1	0																					
DOCUMENT #2	0	0	0																					
DOCUMENT #3	0																							
.																								
.																								
.																								
.																								
DOCUMENT #N																								
ENTRY ID	0	1	2	...	20	...	40	...	60	...	70	...	80	.....	450	...	511							
NUMBER	HIRAGANA CHARACTER HASHING AREA				KATAKANA CHARACTER HASHING AREA				ALPHABET SYMBOL HASHING AREA				NUMERIC CHARACTER HASHING AREA				1ST LEVEL JIS KANJI CHARACTER HASHING AREA				2ND LEVEL JIS KANJI CHARACTER HASHING AREA			



## FIG. 16



## FIG. 17

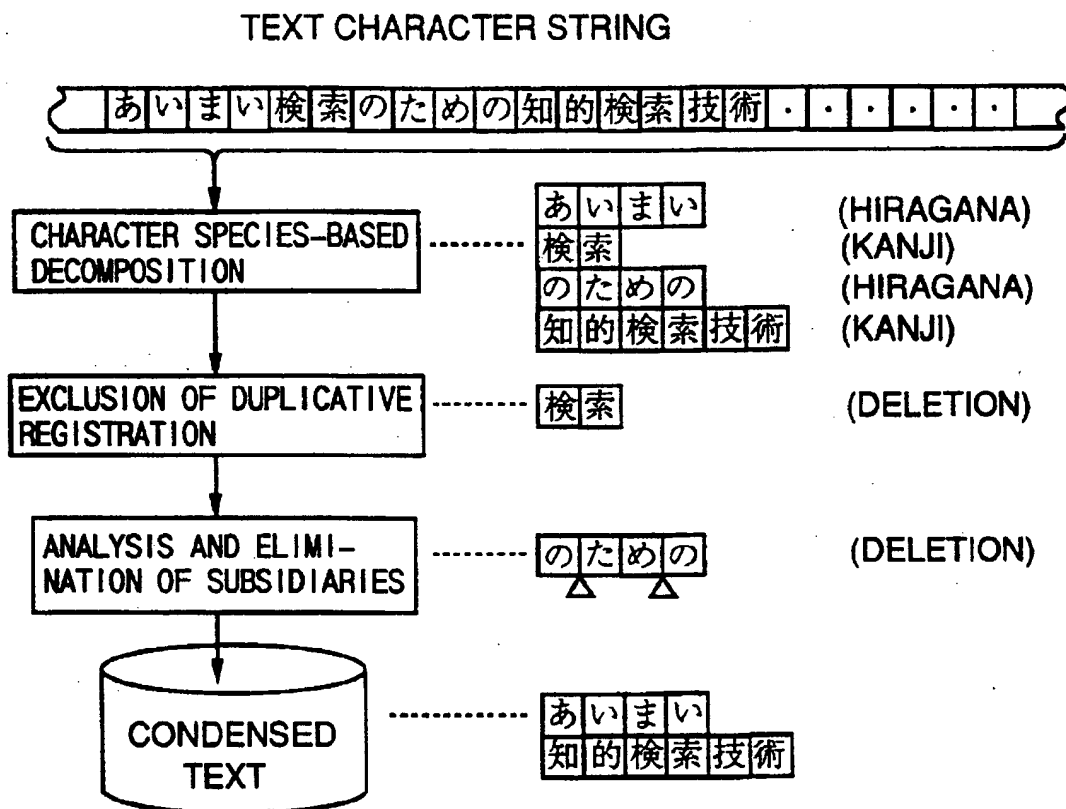


## FIG. 18

検 . . .

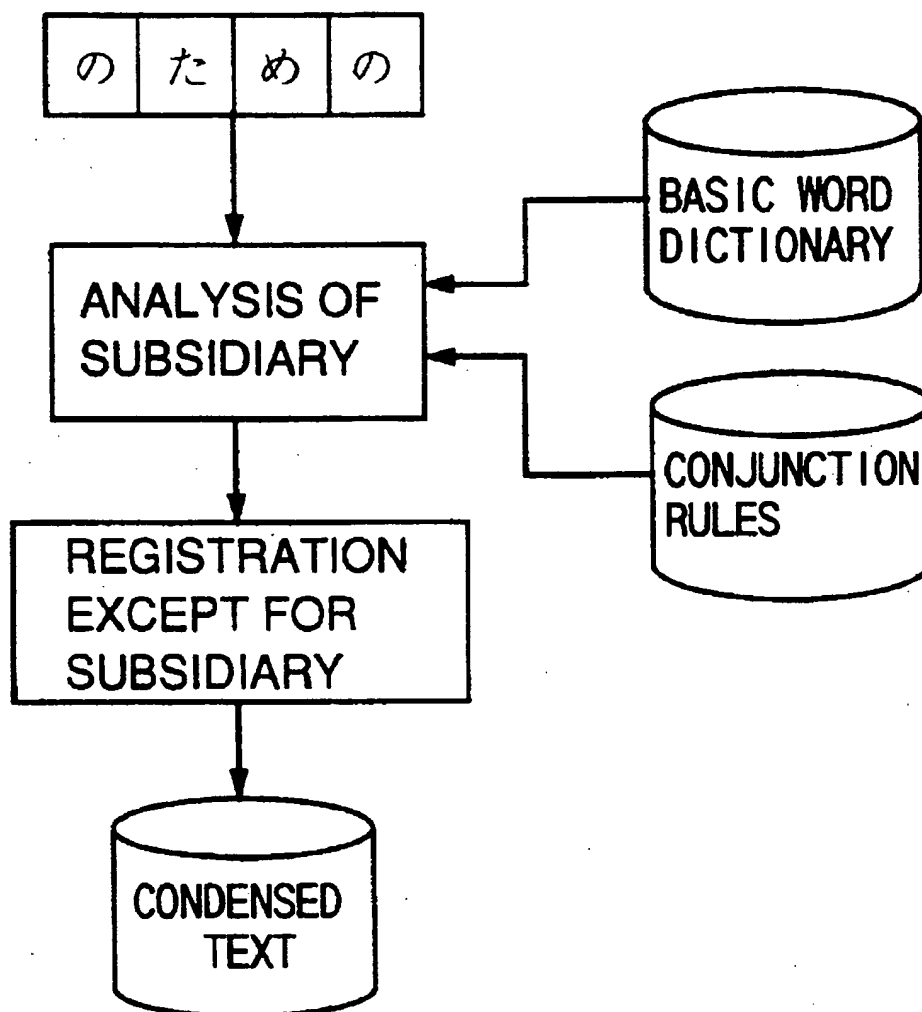
SCODE	ENTRY ID NUMBER
.	.
.	.
.	.
( 095D ) H	150
( 095E ) H	356
( 095F ) H	231
( 0960 ) H	483
( 0961 ) H	2
( 0962 ) H	256
( 0963 ) H	25
( 0964 ) H	67
.	.
.	.
.	.

## FIG. 19



## FIG. 20

## HIRAGANA CHARACTER STRING



## FIG. 21

## BASIC WORD DICTIONARY

〈 PARTICIPIAL FORM OF VERB “ある” 〉	: ある
〈 SUBJUNCTIVE FORM OF VERB “ある” 〉	: あれ
〈 UNENDED FORM OF VERB “なる” 〉	: なら
〈 CONTINUATIVE FORM OF VERB “なる” 〉	: なり
〈 UNENDED FORM OF VERB “もつ” 〉	: もた
〈 POSTPOSITIONAL WORD “が” 〉	: が
〈 NOUN “こと” 〉	: こと
〈 NOUN “ため” 〉	: ため
〈 NOUN “の” 〉	: の

## FIG. 22

## CONJUNCTION RULES

- CONJUNCTION RULE 1 : < PARTICIPIAL FORM OF VERB "ある" > + < NOUN "こと" >  
 CONJUNCTION RULE 2 : < PARTICIPIAL FORM OF VERB "もつ" > + < NOUN "ため" >  
 CONJUNCTION RULE 3 : < NOUN "こと" > + < POSTPOSITIONAL WORD "が" >  
 CONJUNCTION RULE 5 : < PARTICIPIAL FORM OF VERB "する" > + < NOUN "こと" >  
 CONJUNCTION RULE 6 : < POSTPOSITIONAL WORD "の" > + < NOUN "ため" >  
 CONJUNCTION RULE 7 : < NOUN "ため" > + < POSTPOSITIONAL WORD "の" >  
 CONJUNCTION RULE 8 : < CONJUNCTION > + < POSTPOSITIONAL WORD "は" >  
 CONJUNCTION RULE 9 : < UNENDED FORM OF VERB "する" > + < POSTPOSITIONAL WORD "は" >  
 CONJUNCTION RULE 10 : < NOUN "こと" > + < POSTPOSITIONAL WORD "に" >  
 CONJUNCTION RULE 11 : < POSTPOSITIONAL WORD "で" > + < ENDED FORM OF VERB "ある" >  
 CONJUNCTION RULE 12 : < UNENDED FORM OF VERB "ある" > + < POSTPOSITIONAL WORD "は" >  
 CONJUNCTION RULE 13 : < POSTPOSITIONAL WORD "ある" > + < ENDED FORM OF VERB "だ" >

## FIG. 23

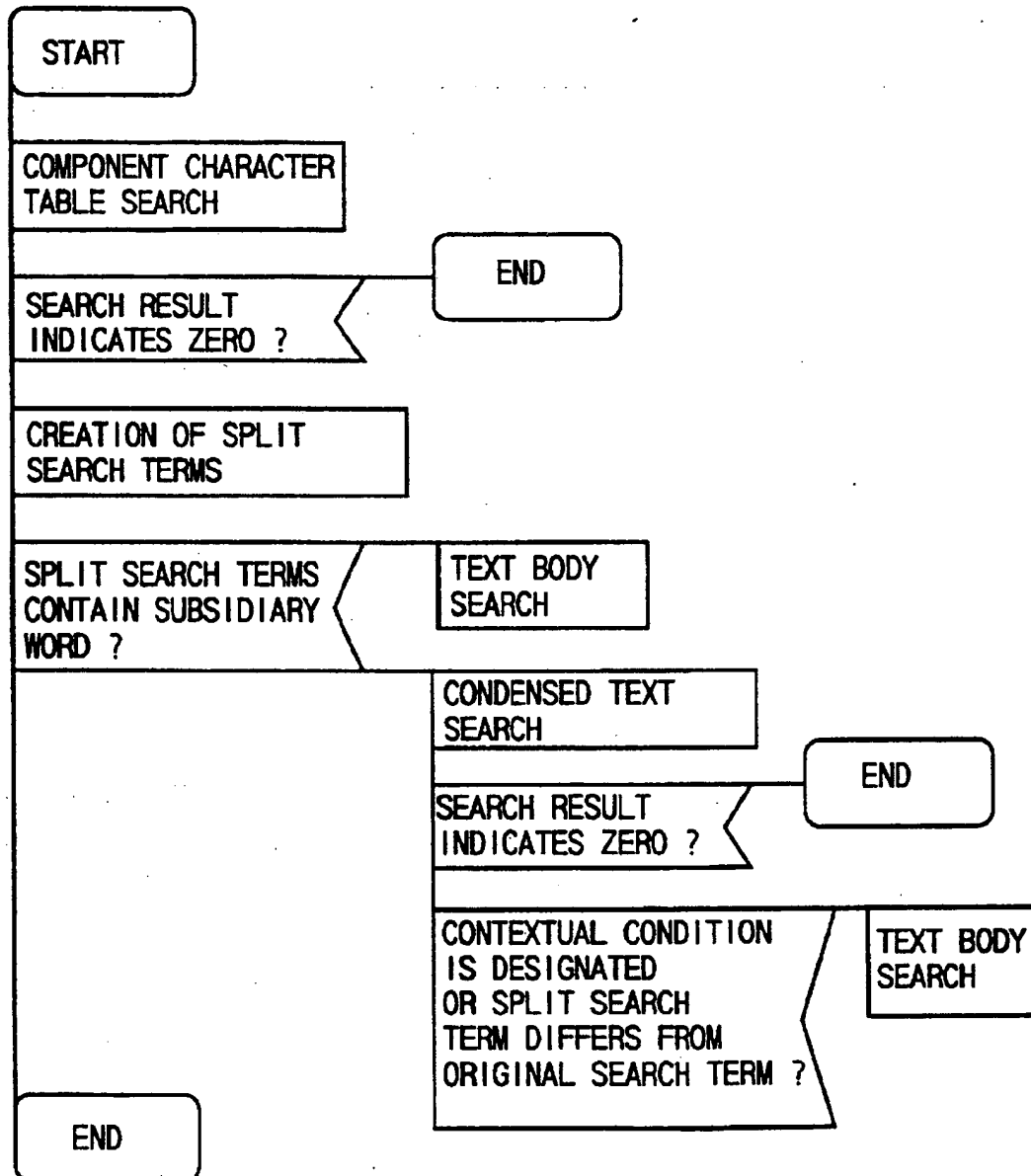
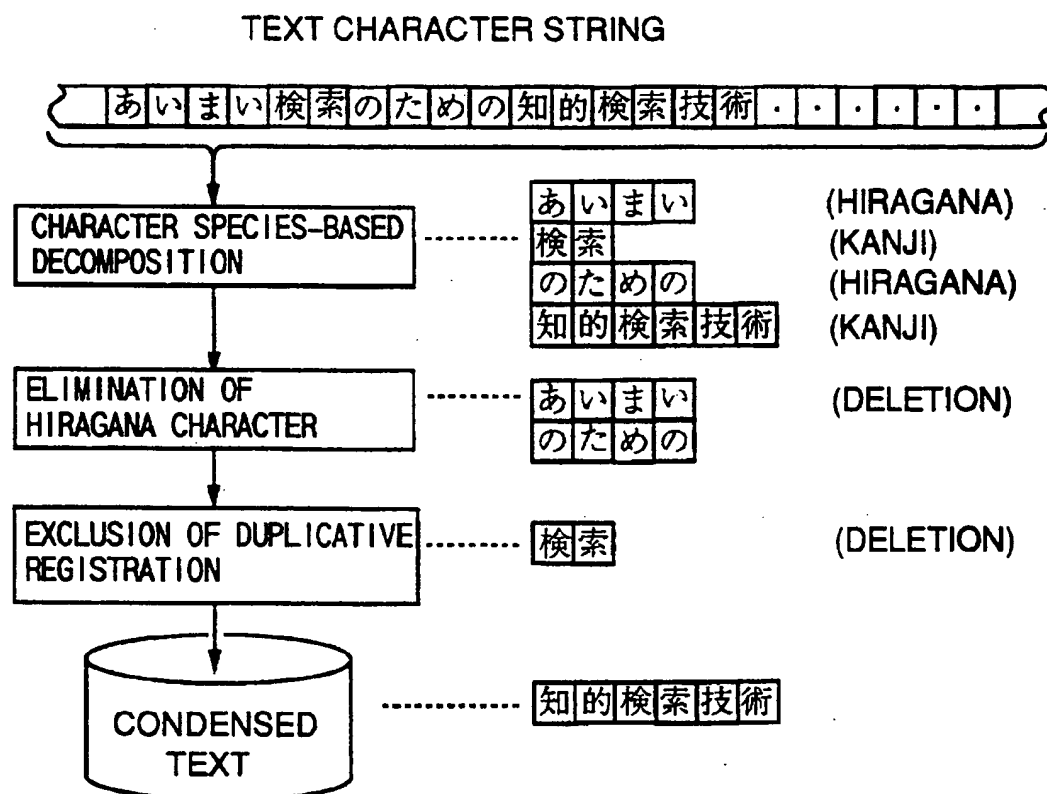




FIG. 24



## FIG. 25

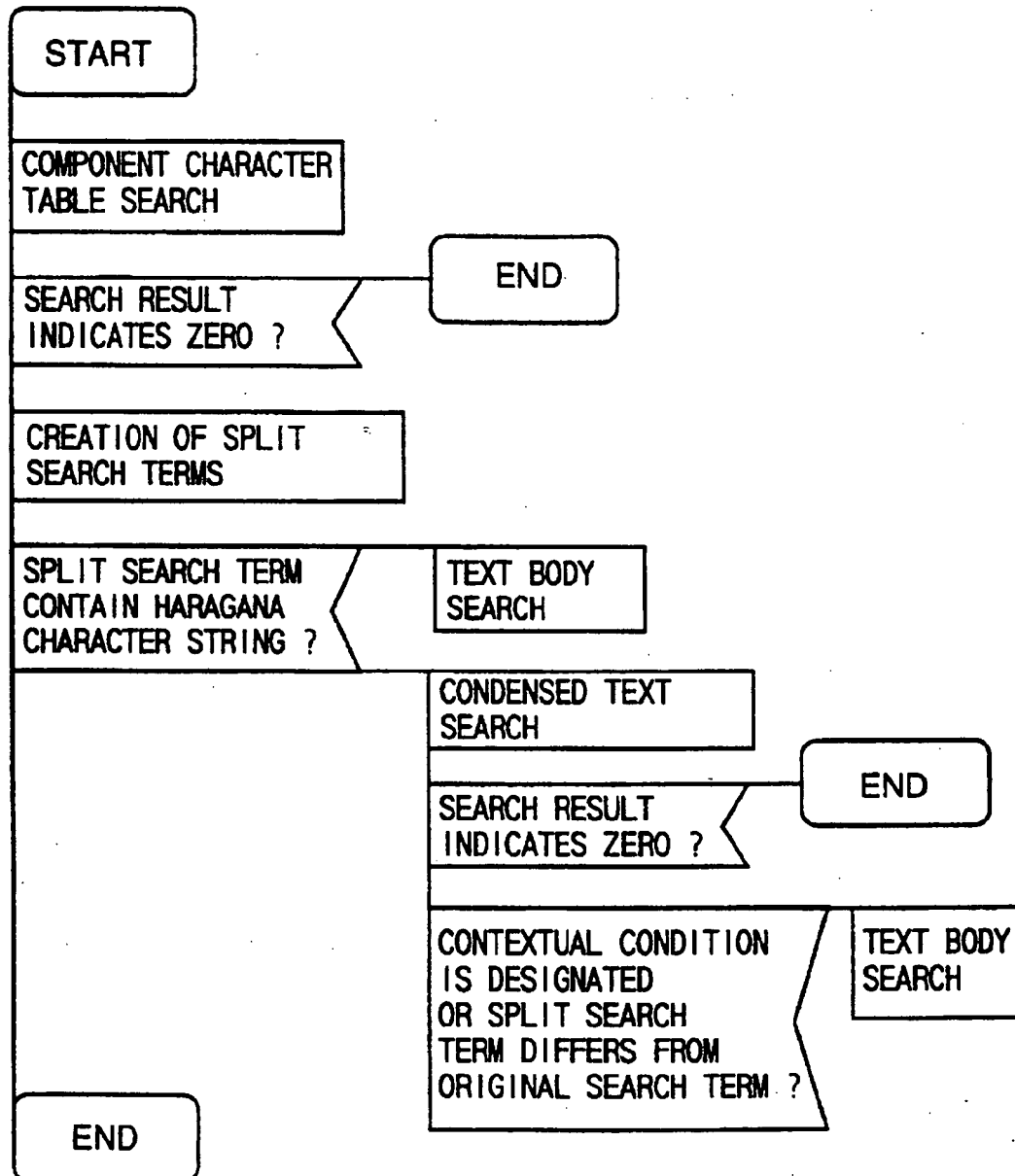
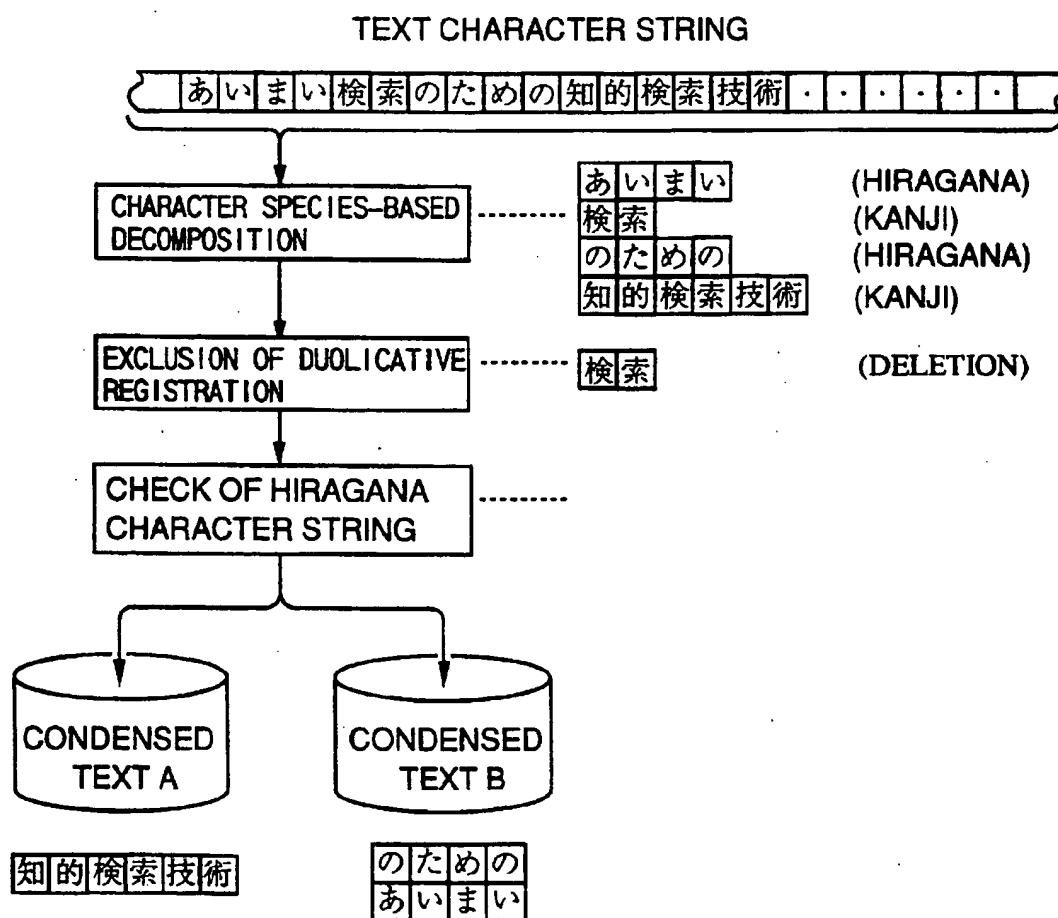


FIG. 26



## FIG. 27

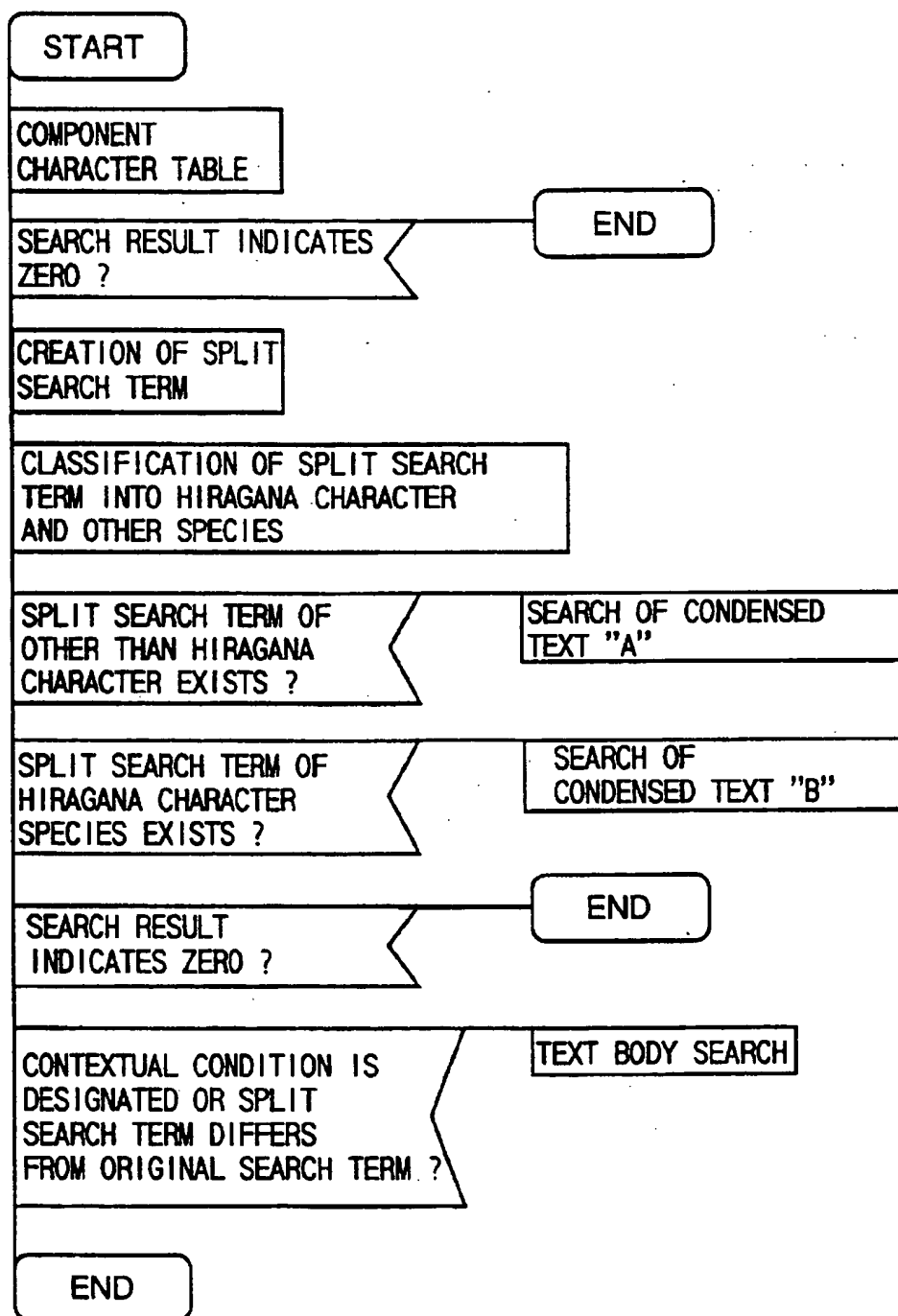
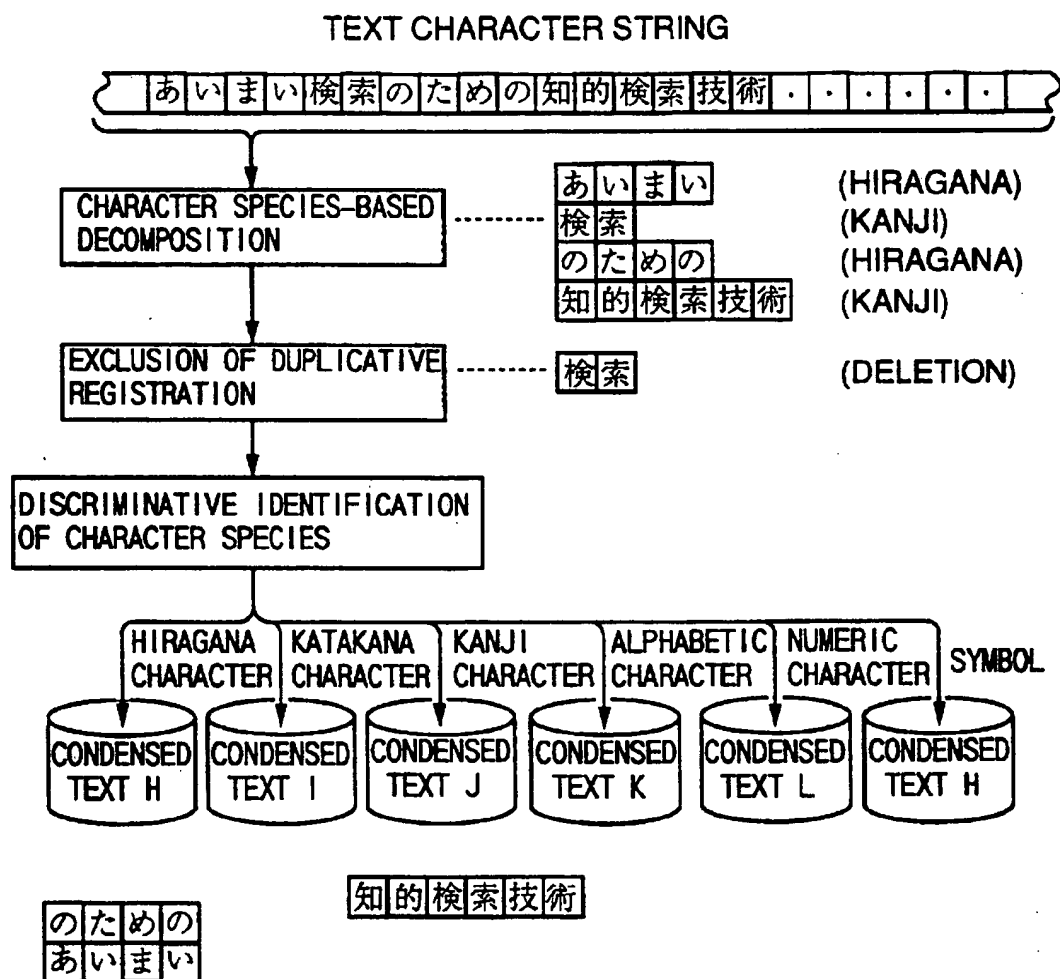


FIG. 28



## FIG. 29

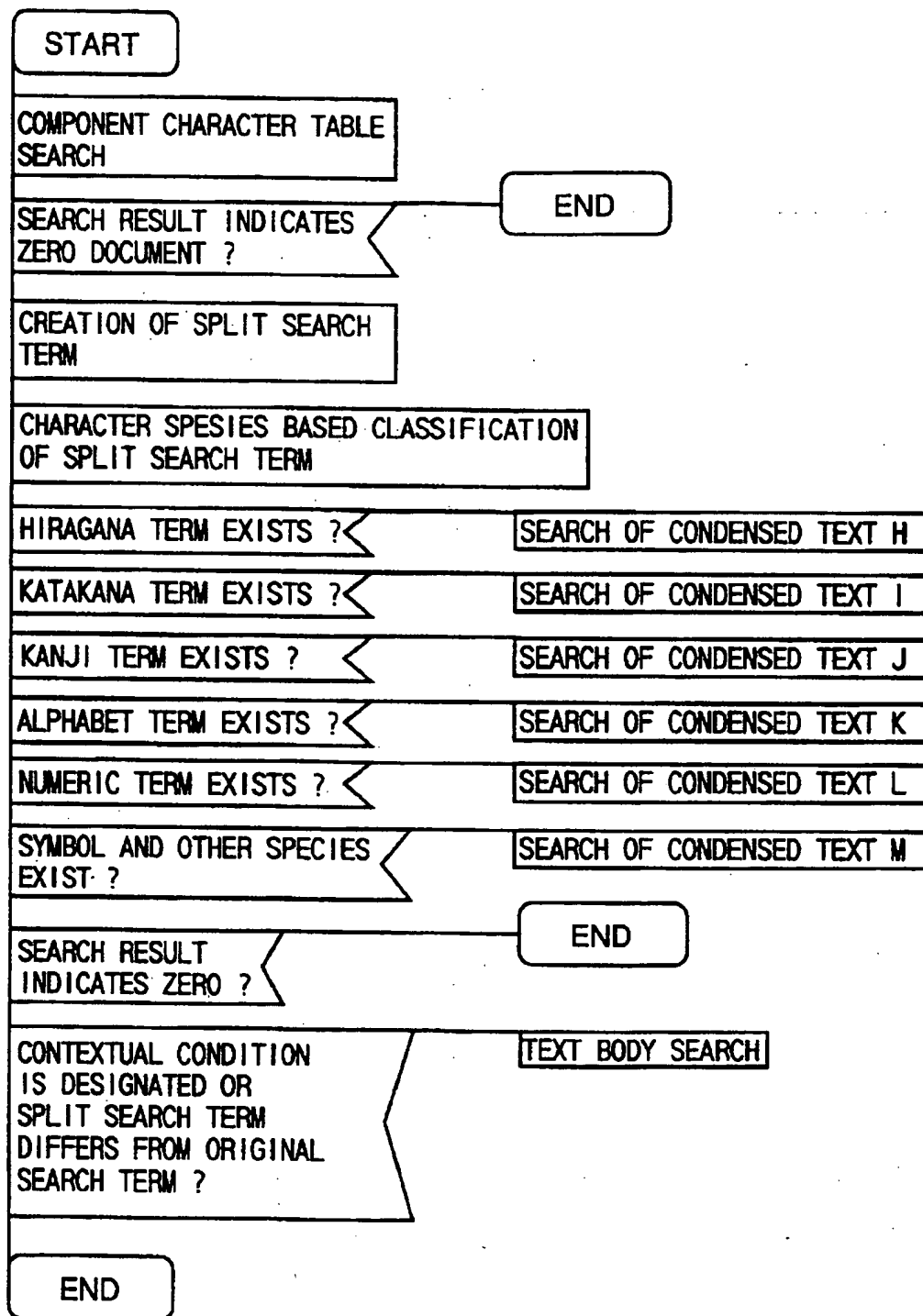


FIG. 30

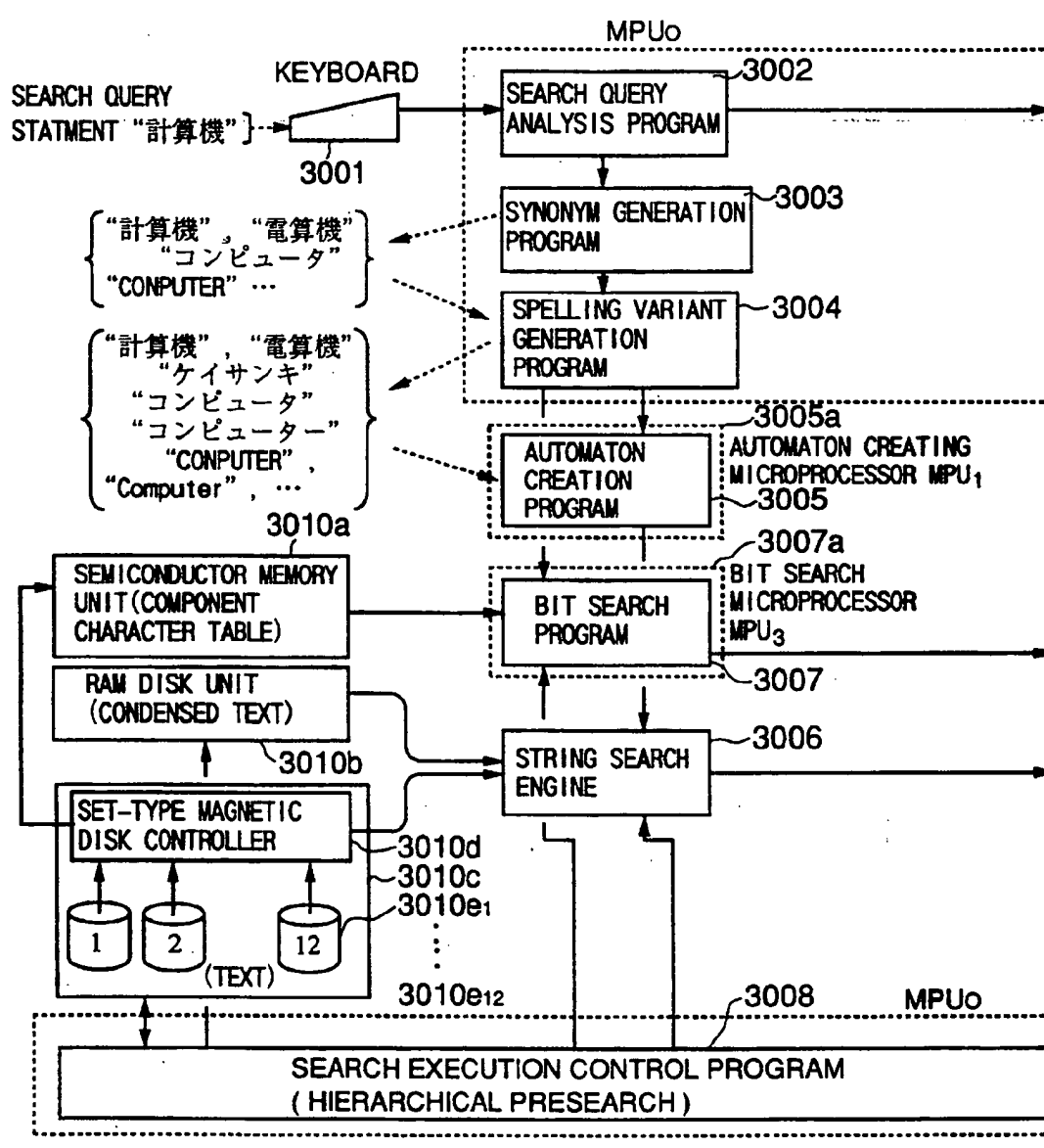


FIG. 31

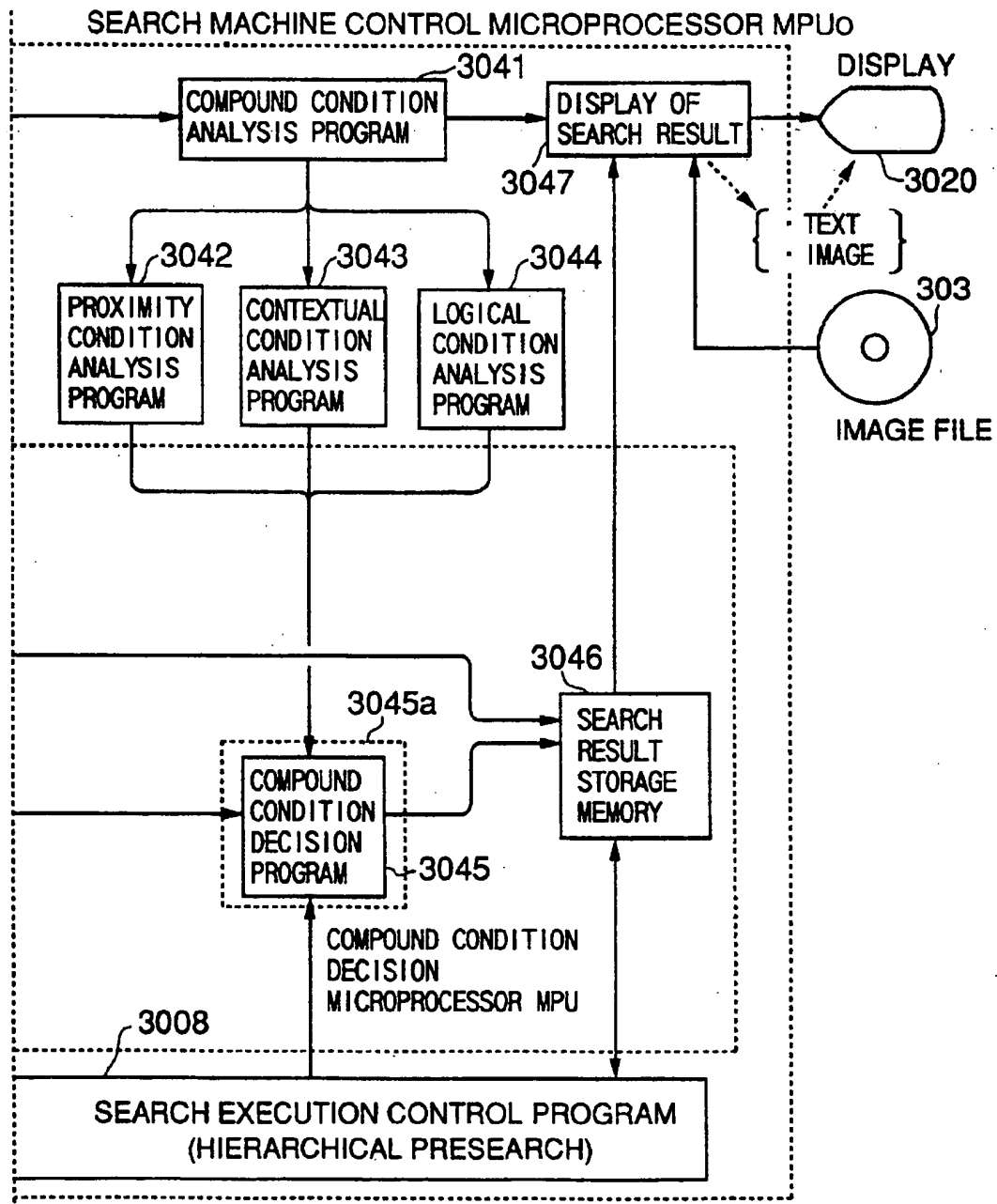
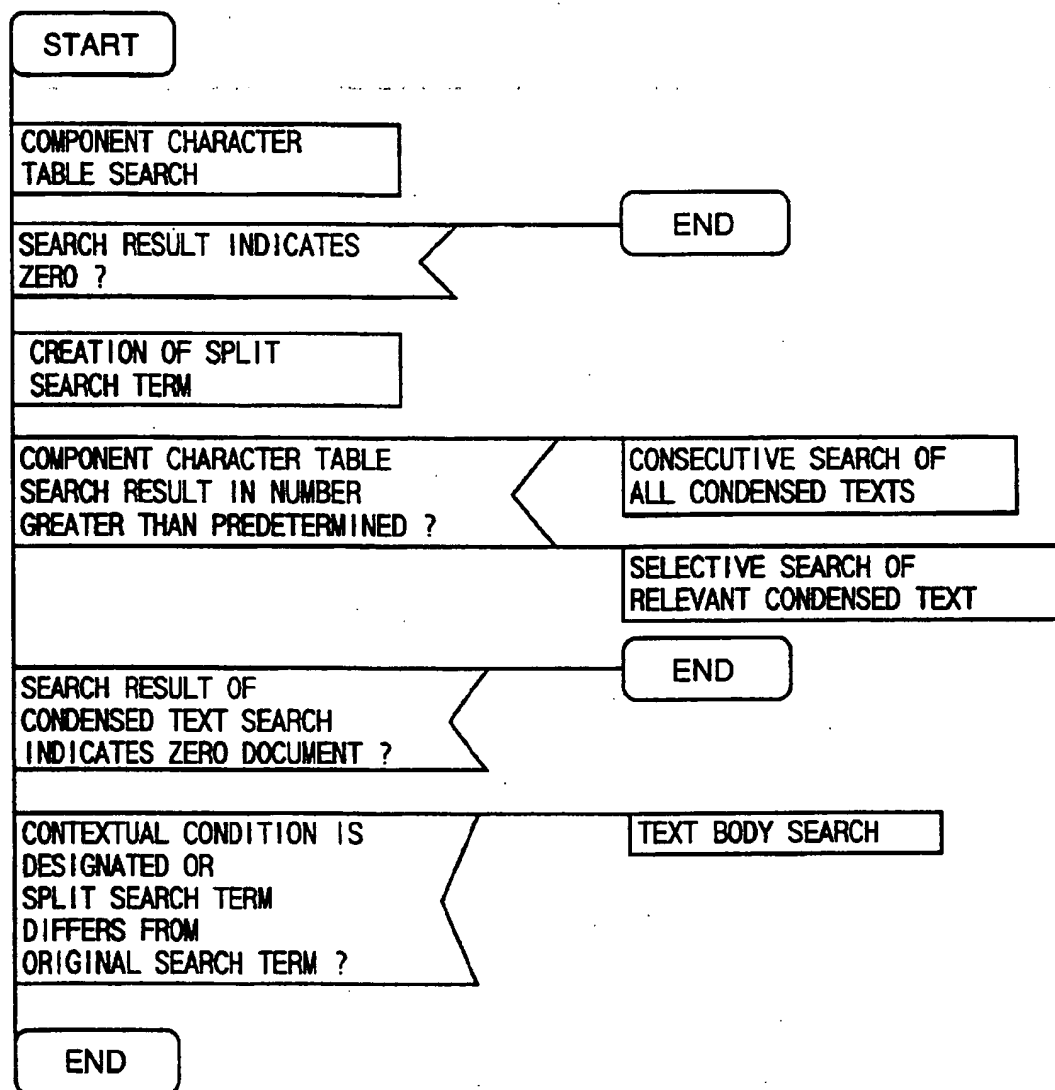




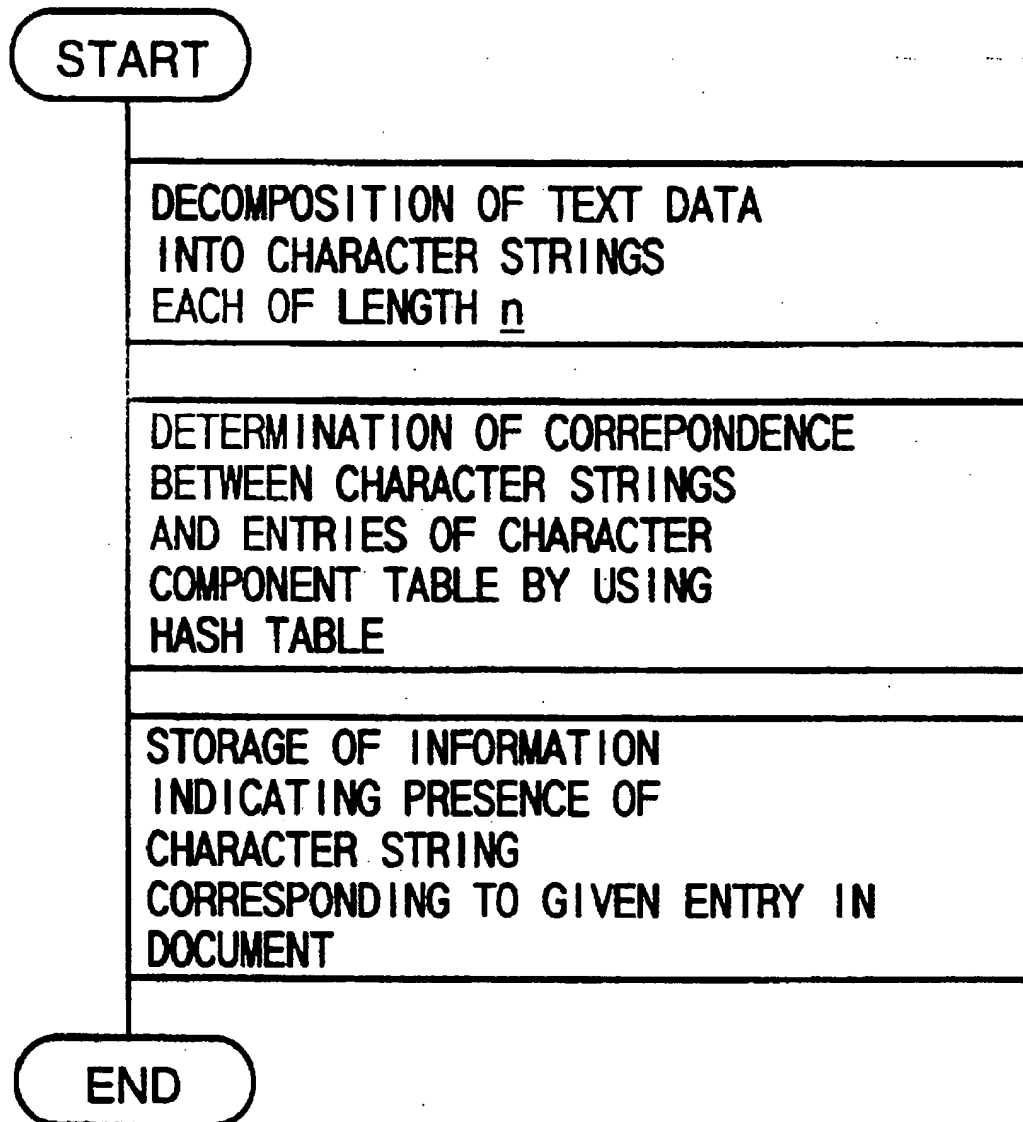
FIG. 32



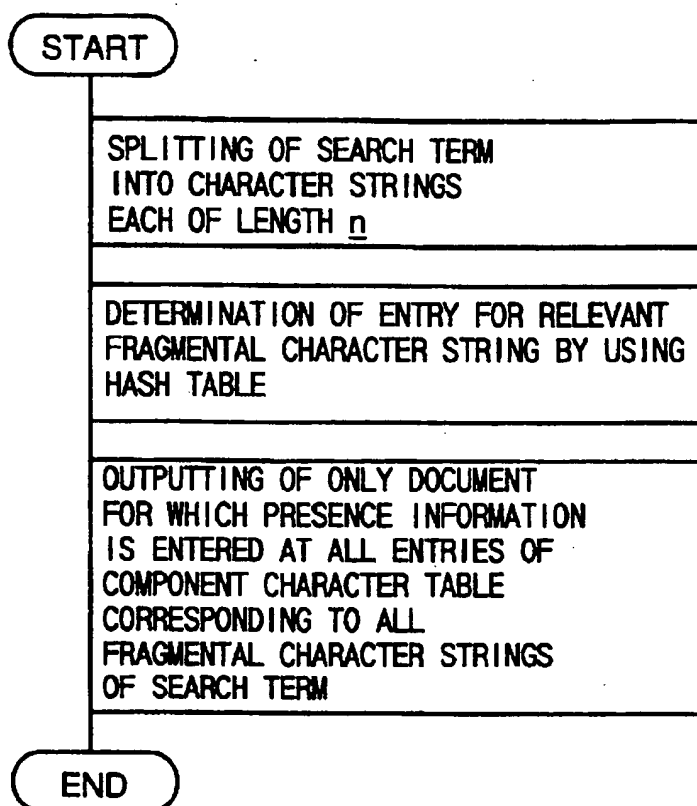
## FIG. 33

DOCUMENT #1	あいま検索のため技術 . . . .
DOCUMENT #2	自然語による検索技術 . . . .
DOCUMENT #3	壁を検出しながら出口,検索 . . . .
DOCUMENT #4	文書理解を用いた検索,システム . . .
.	. . . . .
.	. . .
DOCUMENT #N	

## FIG. 34



## FIG. 35

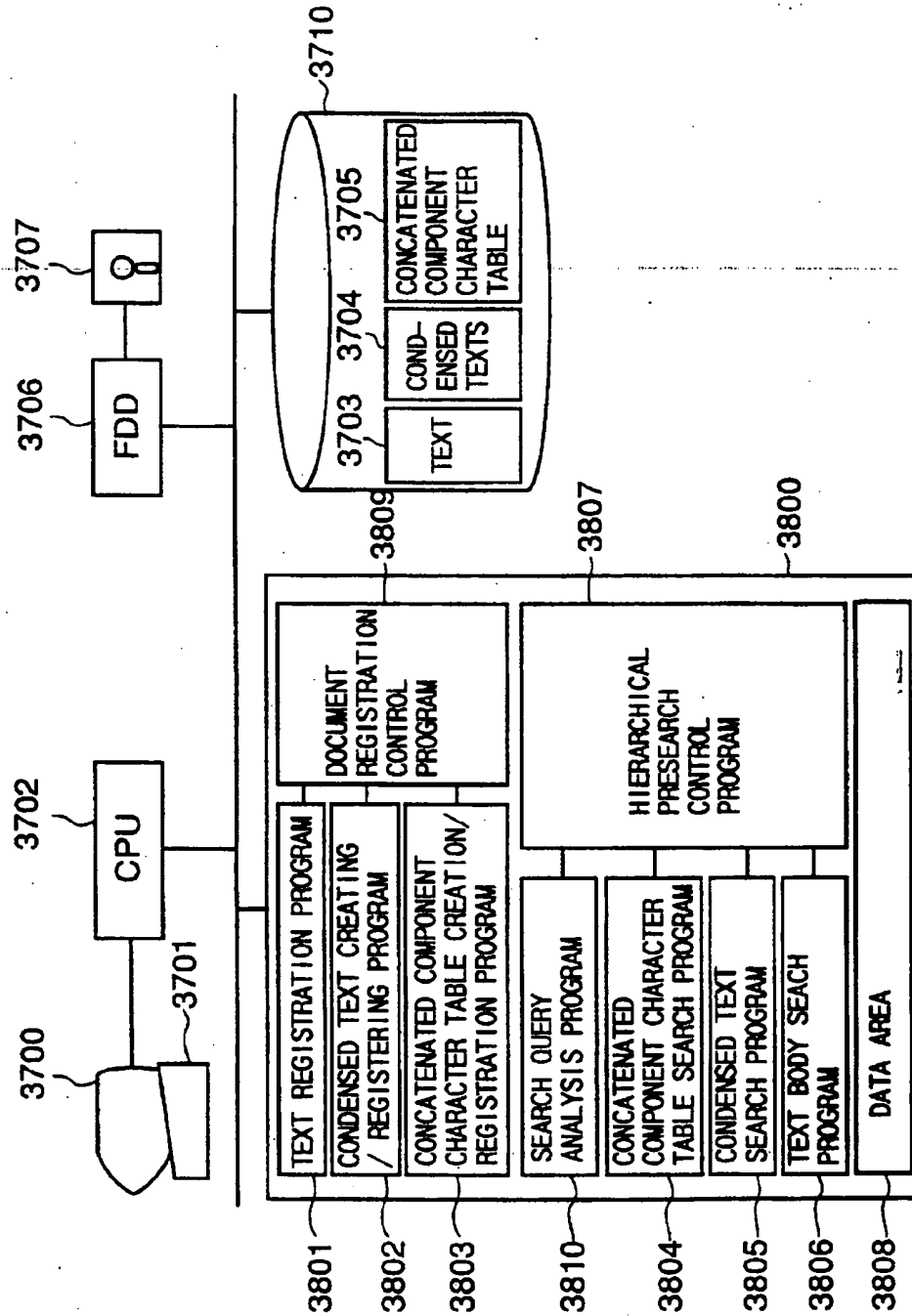


## FIG. 36

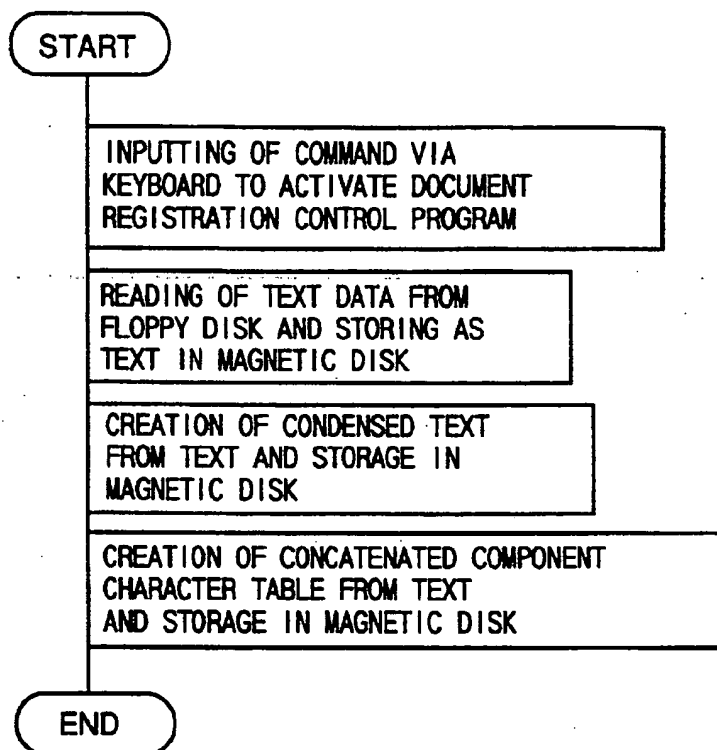
① インターフェース

② インターフェース

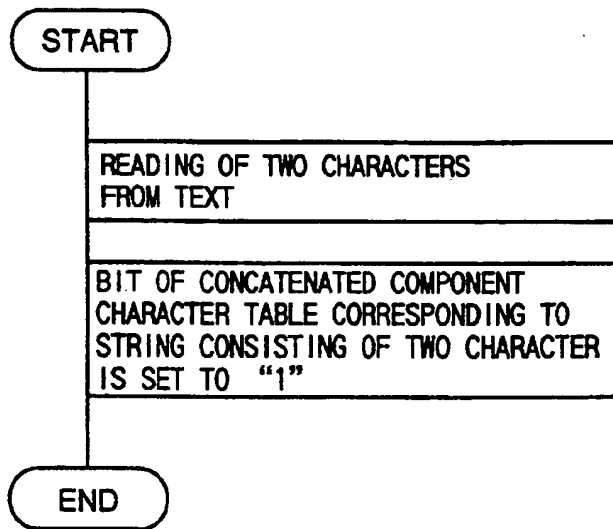
FIG. 37



## FIG. 38



## FIG. 39

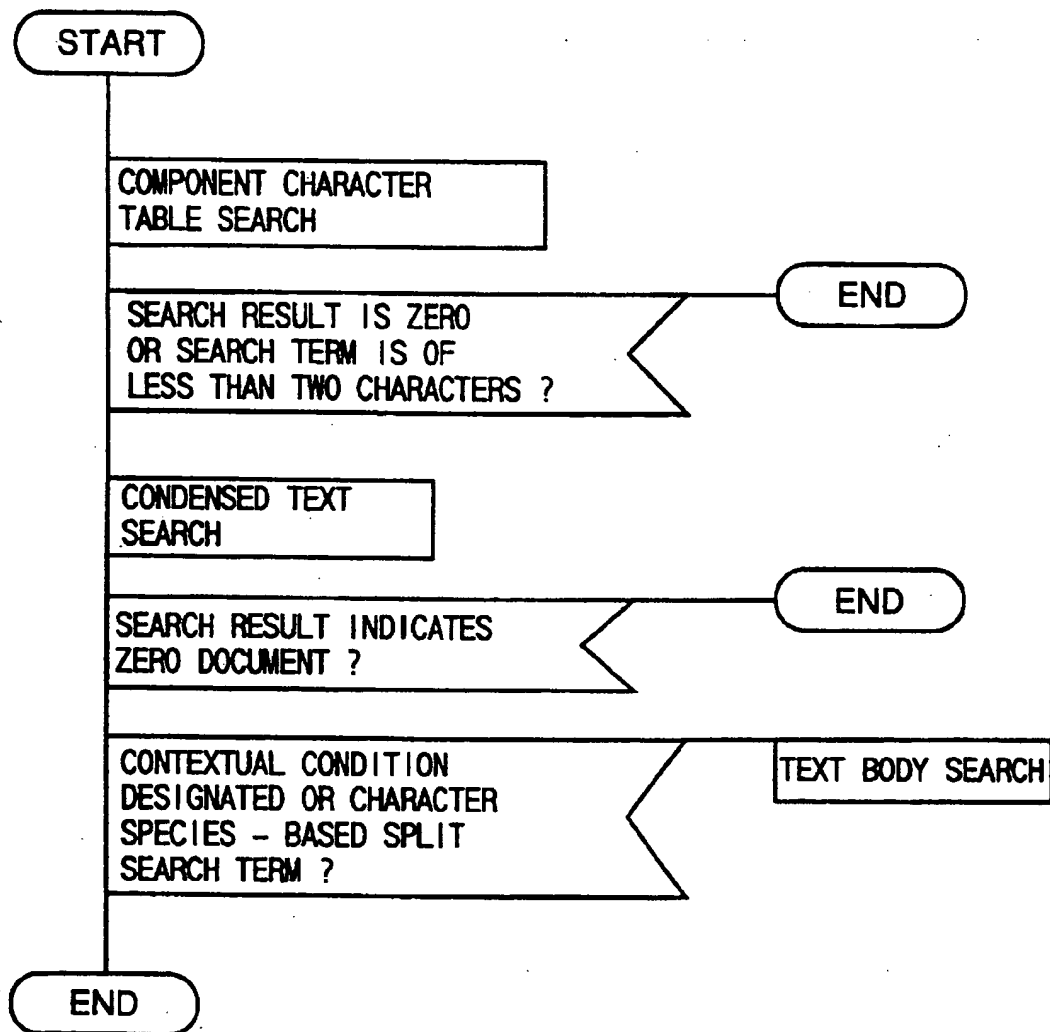


## FIG. 40

DOCUMENT #1	・ ・ オンラインサービスが経済発展の ・ ・
DOCUMENT #2	・ ・ ライオンの生態については未だに ・ ・
DOCUMENT #3	・ ・ 戦争後のイランとイラクの関係は ・ ・

	あああい ・ オン ン イ ラ イ ラン ・ 関係経済生態戦争 ・ ・													
DOCUMENT #1	0	0		1	0	0	1	1	0		0	1	0	0
DOCUMENT #2	0	0		1	1	0	0	1	0		0	0	1	0
DOCUMENT #3	0	0		0	0	1	0	0	1		1	0	0	1
・														
・														
DOCUMENT #N	0	0		0	0	0	0	0	0		0	0	0	0

## FIG. 41







## FIG. 43

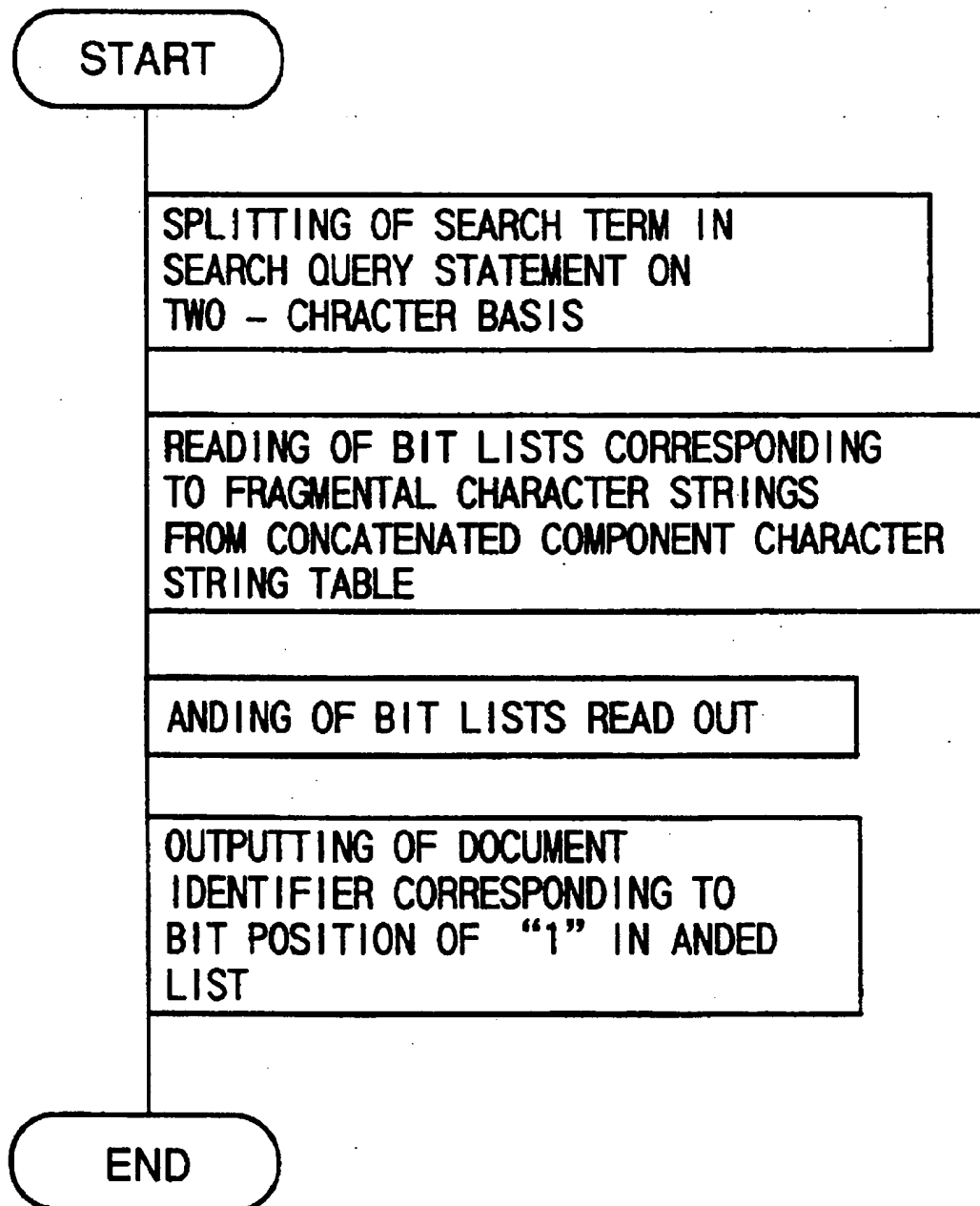


FIG. 44

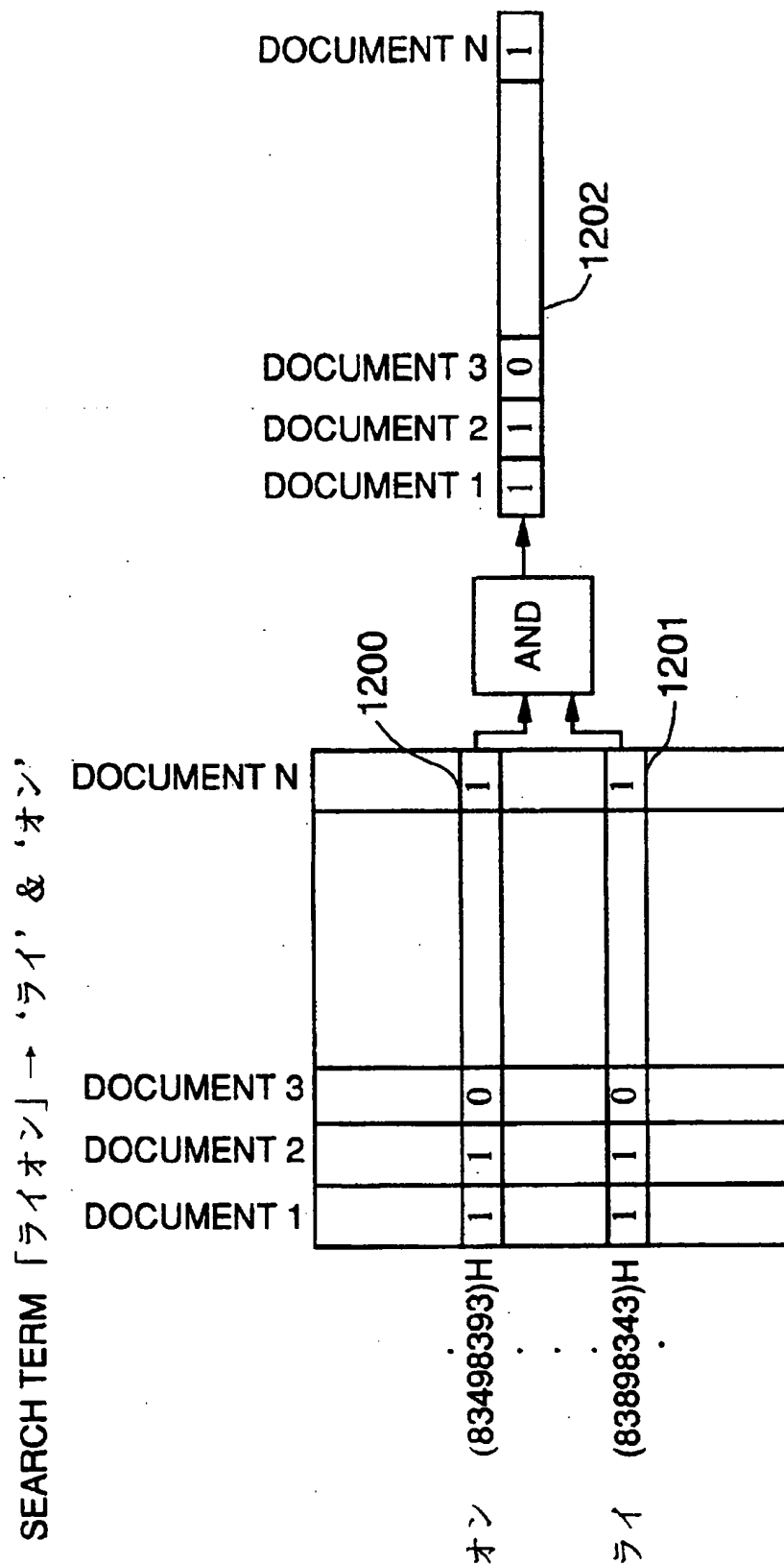


FIG. 45

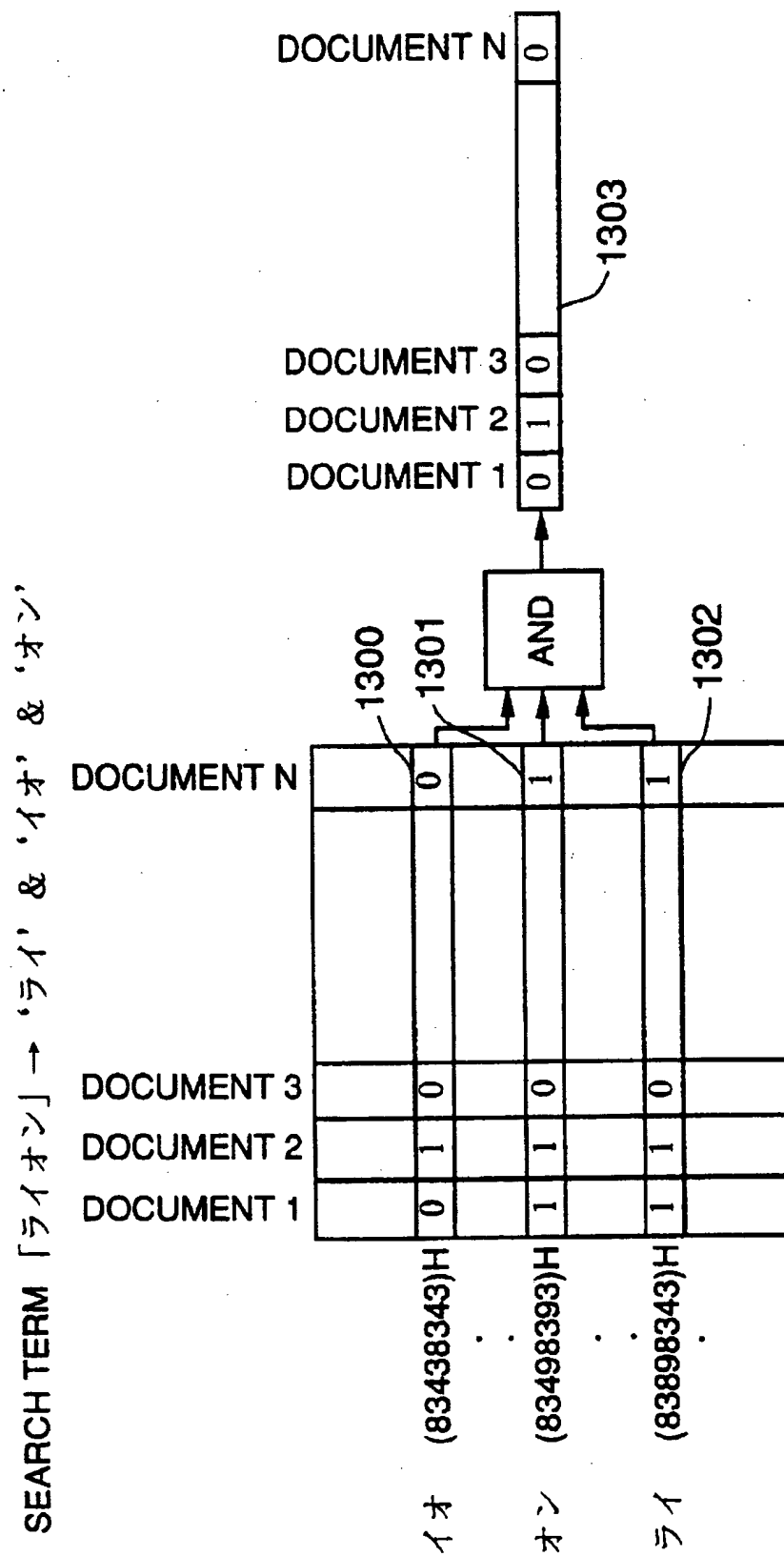


FIG. 46

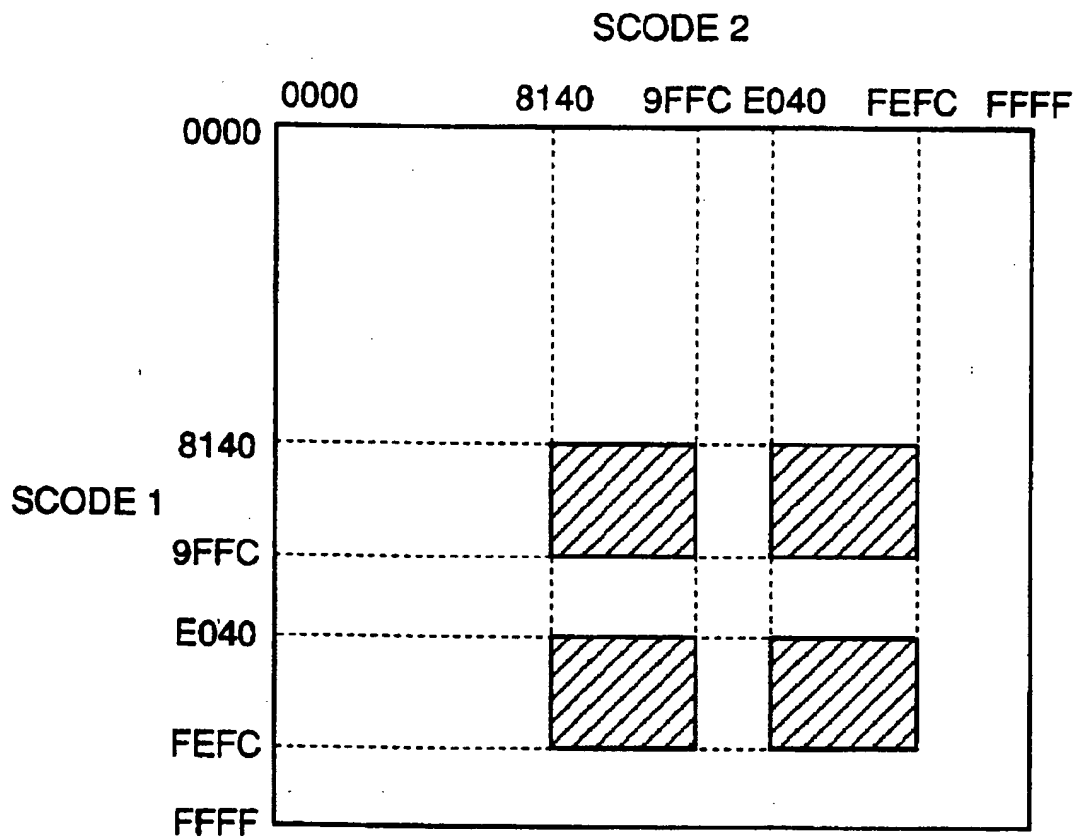
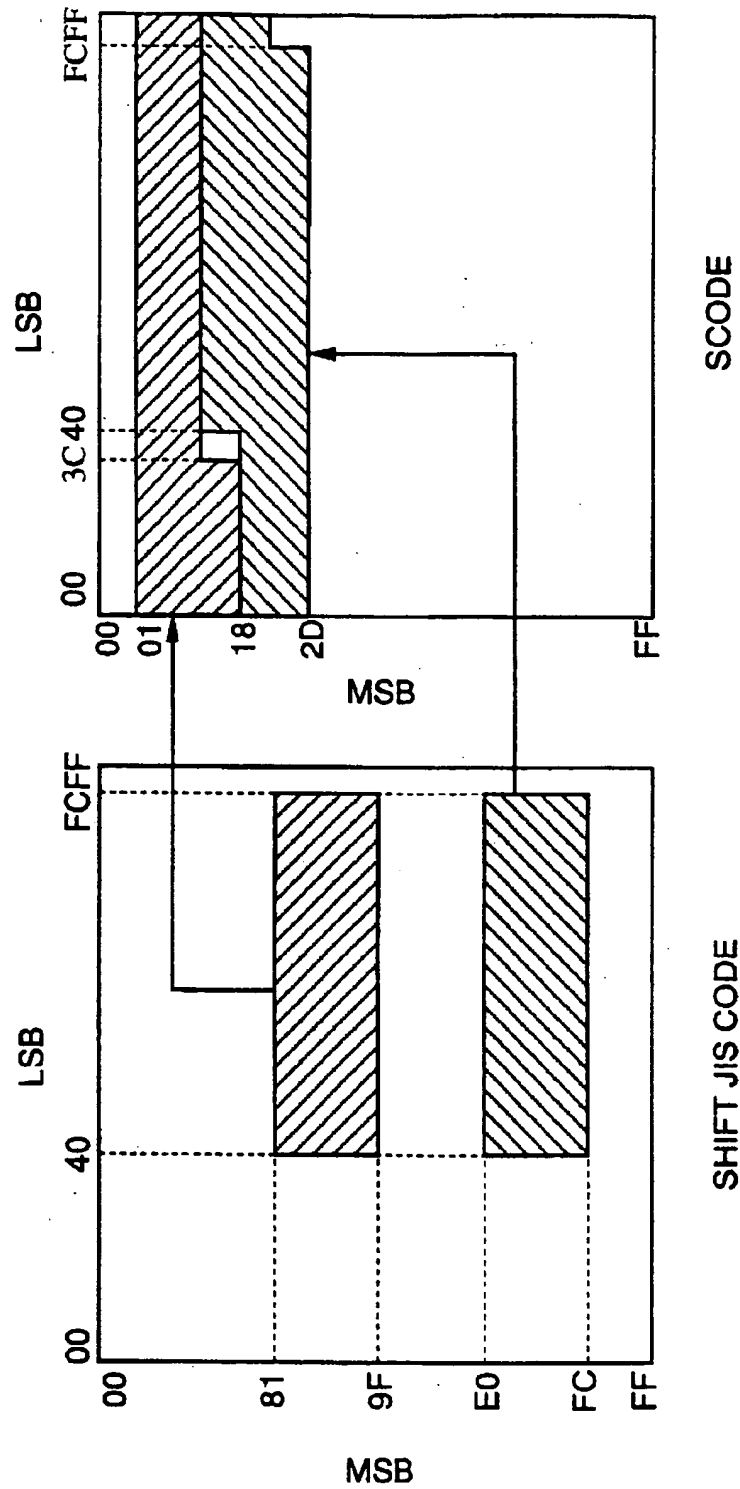


FIG. 47



$$\text{SCODE} = (\text{SJIS\_H} \& 0 \times \text{BF}) * 0 \times \text{C0} + \text{SJIS\_L} - 0 \times 6000$$

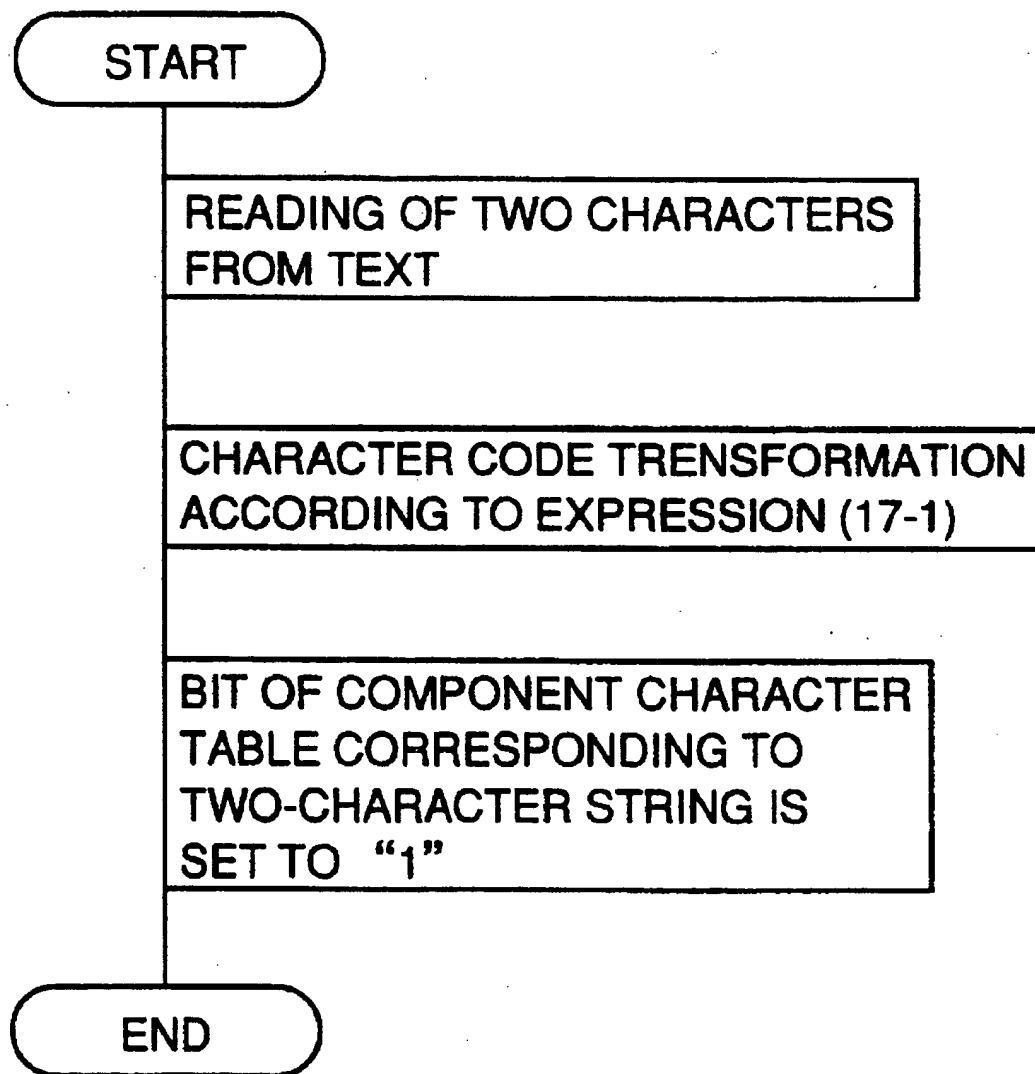
FIG. 48

DOCUMENT #1	・・・オンラインサービスが経済発展の・・・
DOCUMENT #2	・・・ライオンの生息については未だに・・・

COMBINATION OF (00000000)H (028602D3)H (096F0A8F)H (20222022)H  
CHARACTER CODES (02830286)H (02C90283)H

	イ	オ	ン	ヲ	・・・	経済	
DOCUMENT #1	1	1	1	1	・・・	1	0
DOCUMENT #2	0	1	1	1	・・・	0	0
DOCUMENT #3							
...							
DOCUMENT #N	0	0	0	0	・・・	0	0

## FIG. 49





## FIG. 50

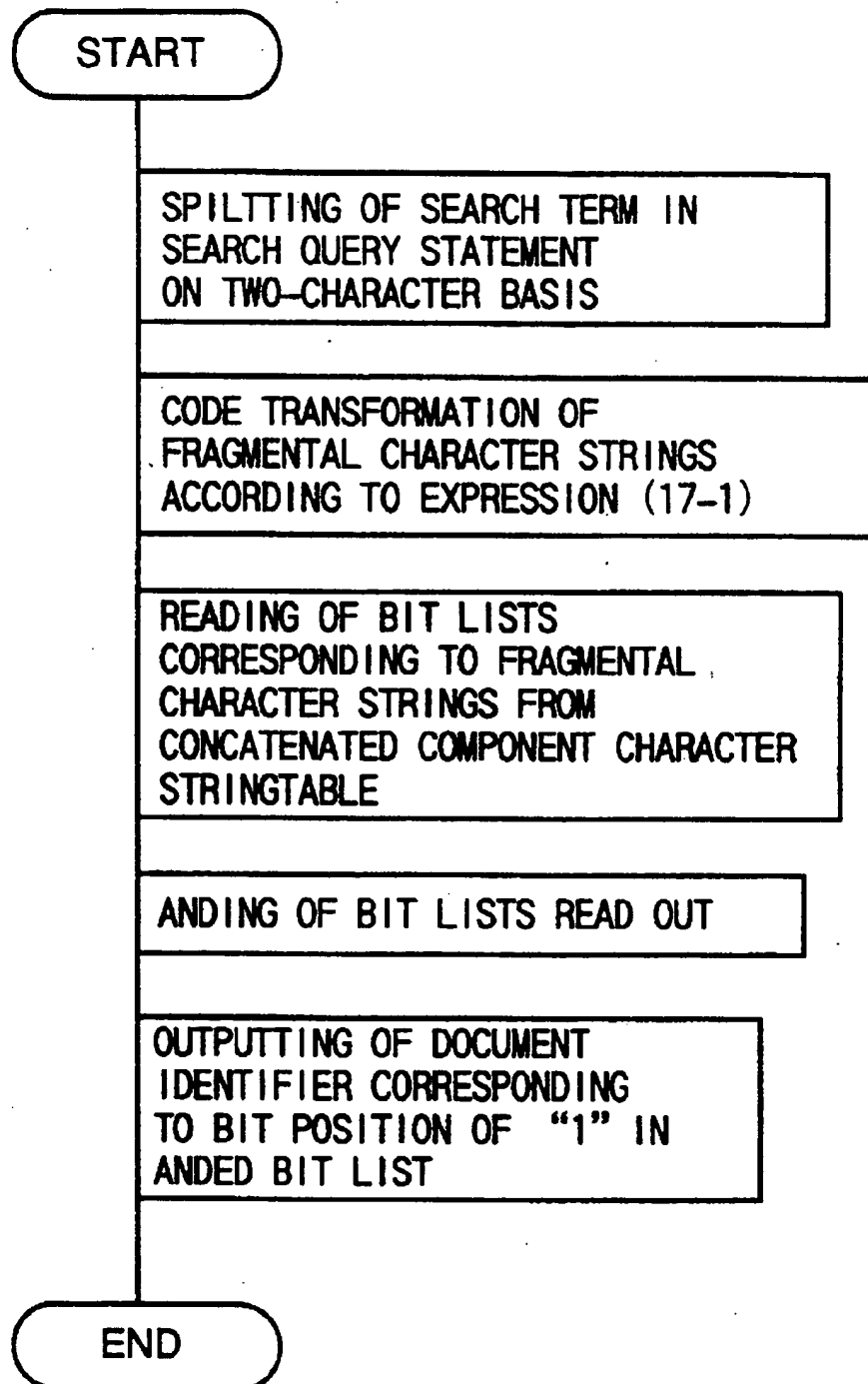
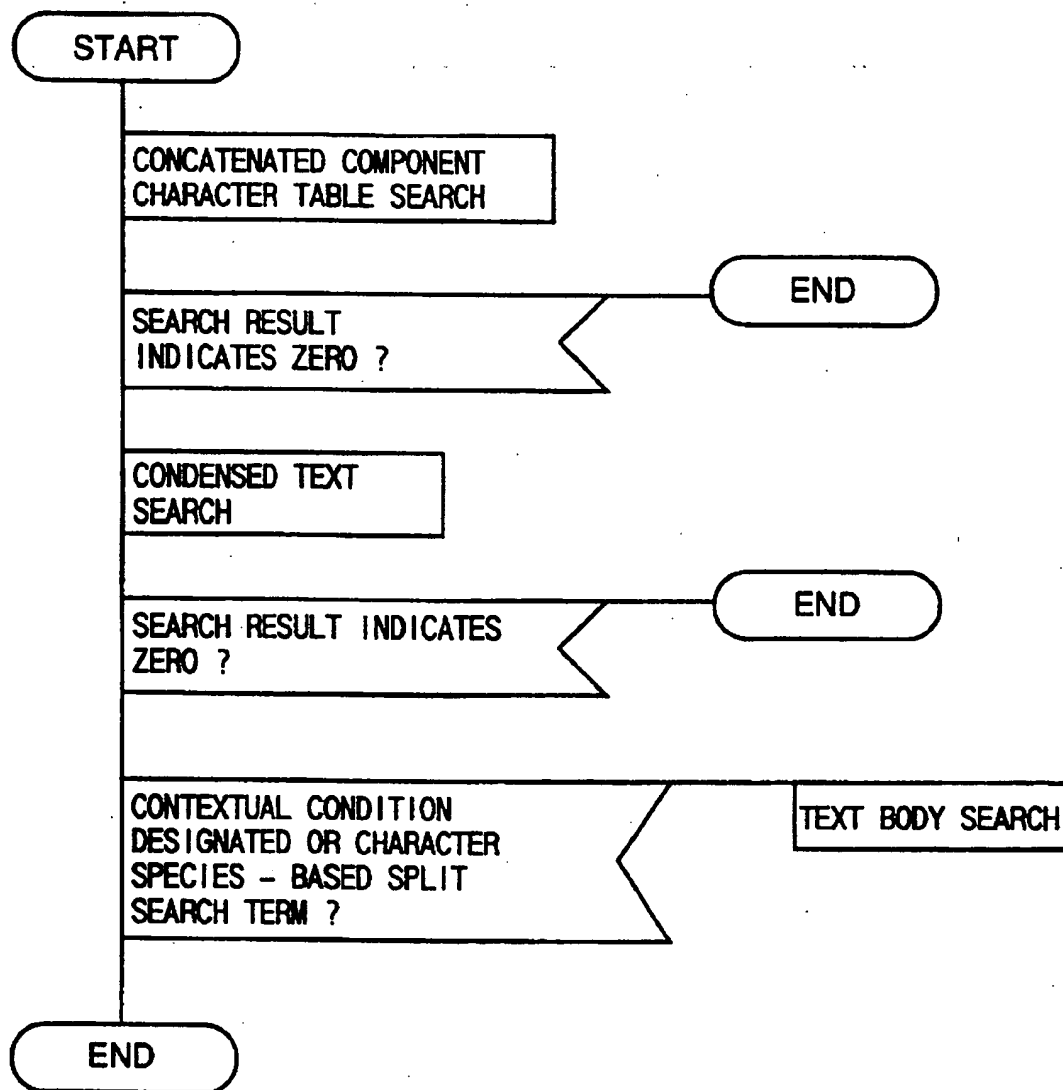
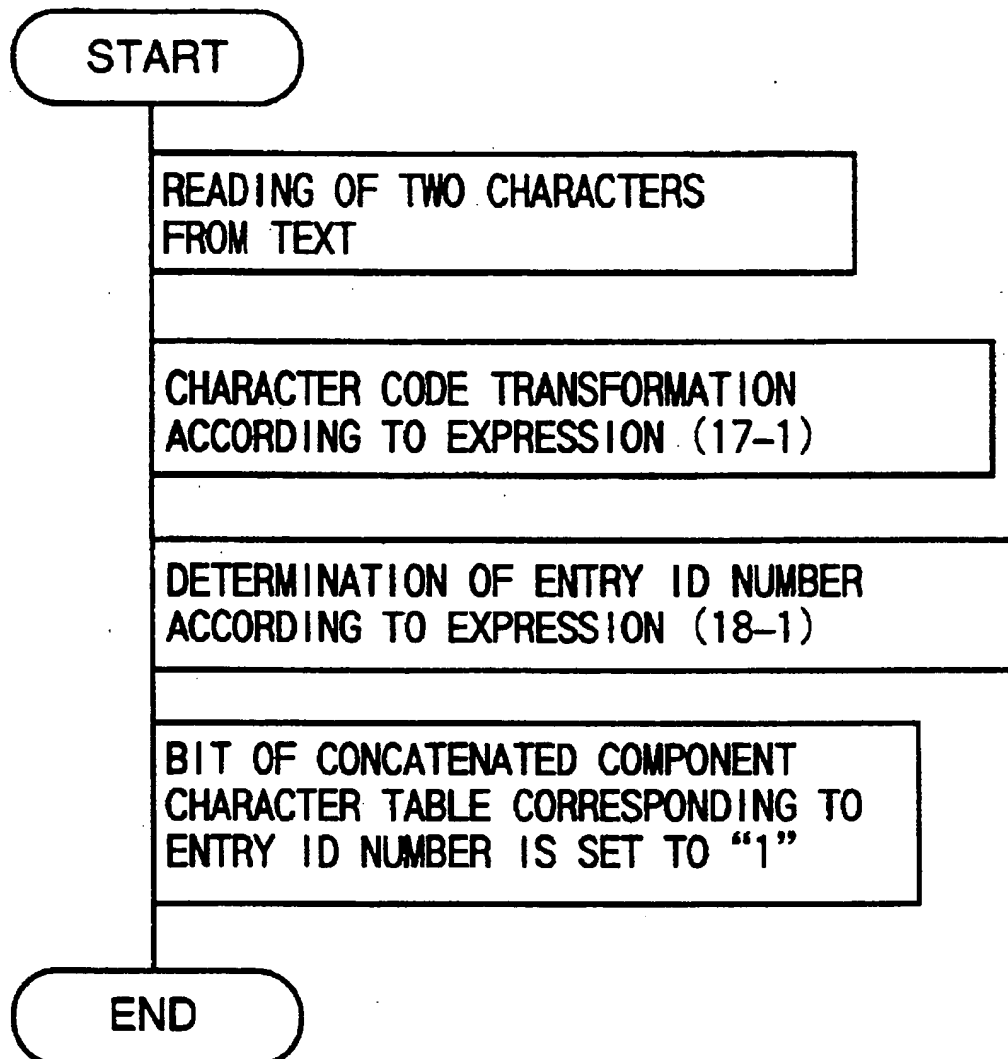


FIG. 51



## FIG. 52



# FIG. 53

DOCUMENT #1	...オンラインサービスが経済発展の...
DOCUMENT #2	...ライオンの生息については未だに...

HASH ENTRY	0	34	643	646	723	2703	4095
DOCUMENT #1	0	1	1	1	1	1	1
DOCUMENT #2	0	0	1	1	1	0	1
DOCUMENT #3							
DOCUMENT #N	0	0	0	0	0	0	1

## FIG. 54

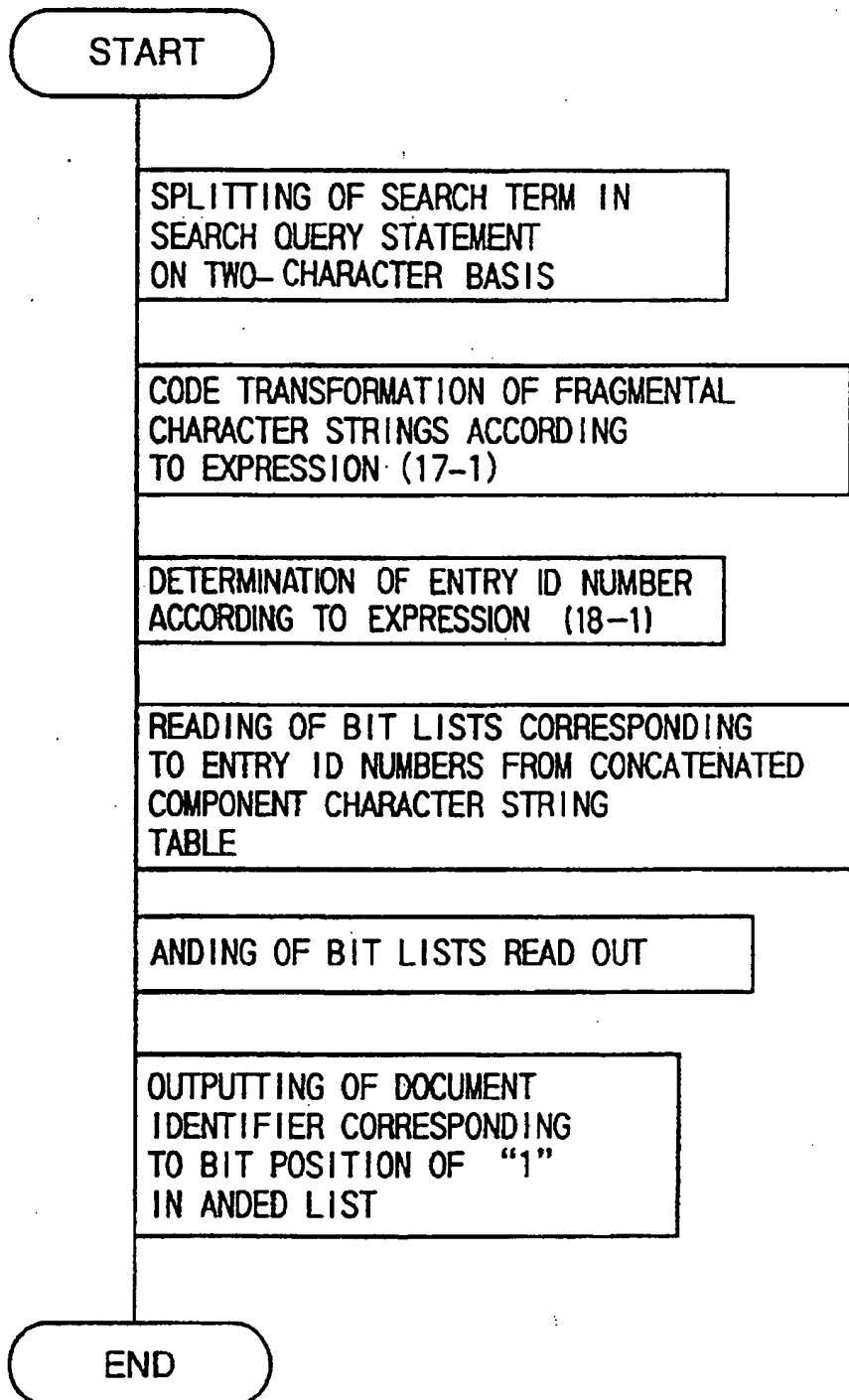


FIG. 55

HASH ENTRY	0	200	400	500	550	2050	2150	4095
DOCUMENT #1								
DOCUMENT #2								
DOCUMENT #3								
DOCUMENT #N								
	HIRAGANA CHARACTER HASHING AREA	KATAKANA CHARACTER HASHING AREA	ALPHABET HASHING AREA	NUMERIC CHARACTER HASHING AREA	1ST LEVEL JIS KANJI CHARACTER HASHING AREA	2ND LEVEL JIS KANJI CHARACTER HASHING AREA	COMPOUND CHARACTER SPECIES HASHONG AREA	

## FIG. 56

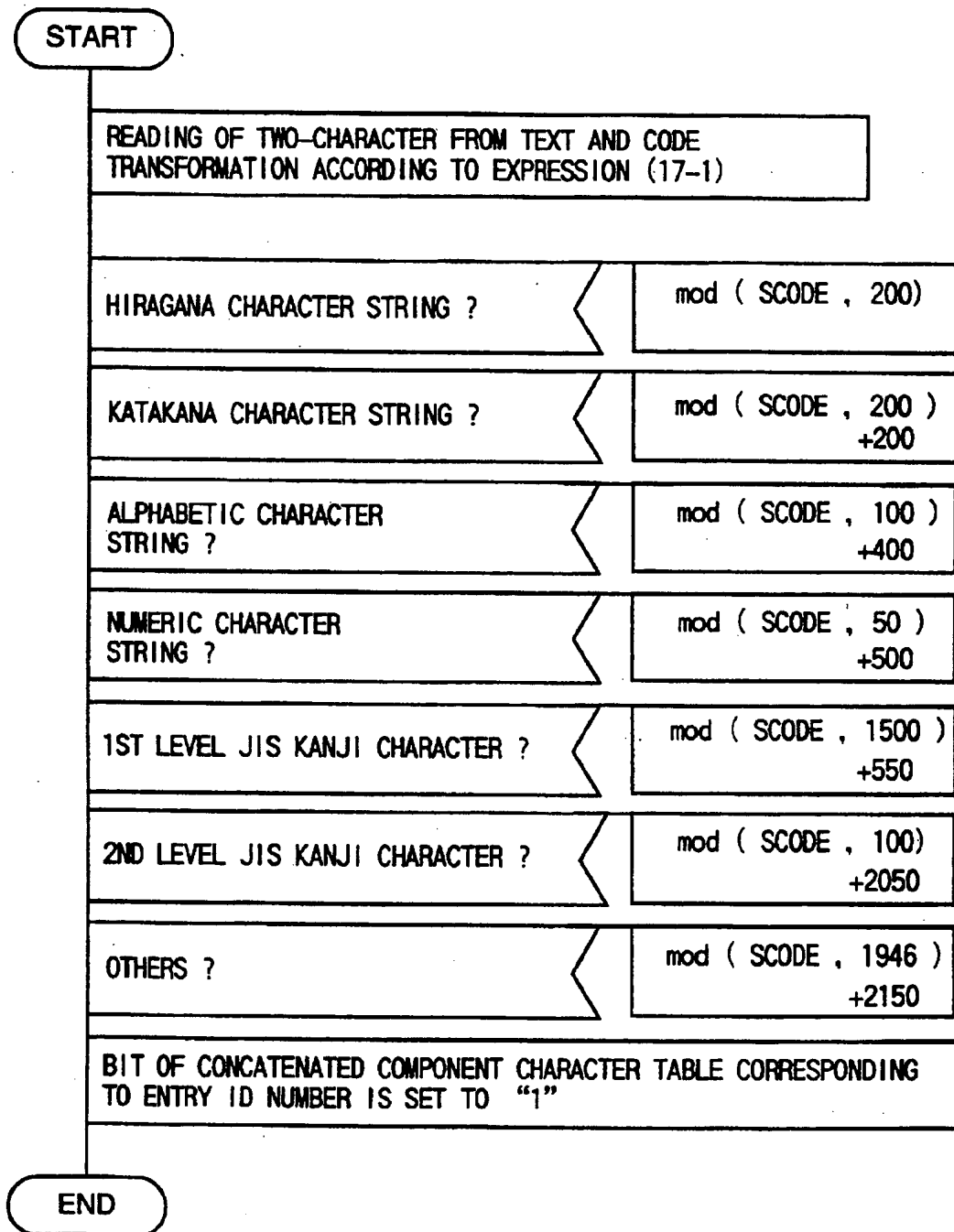
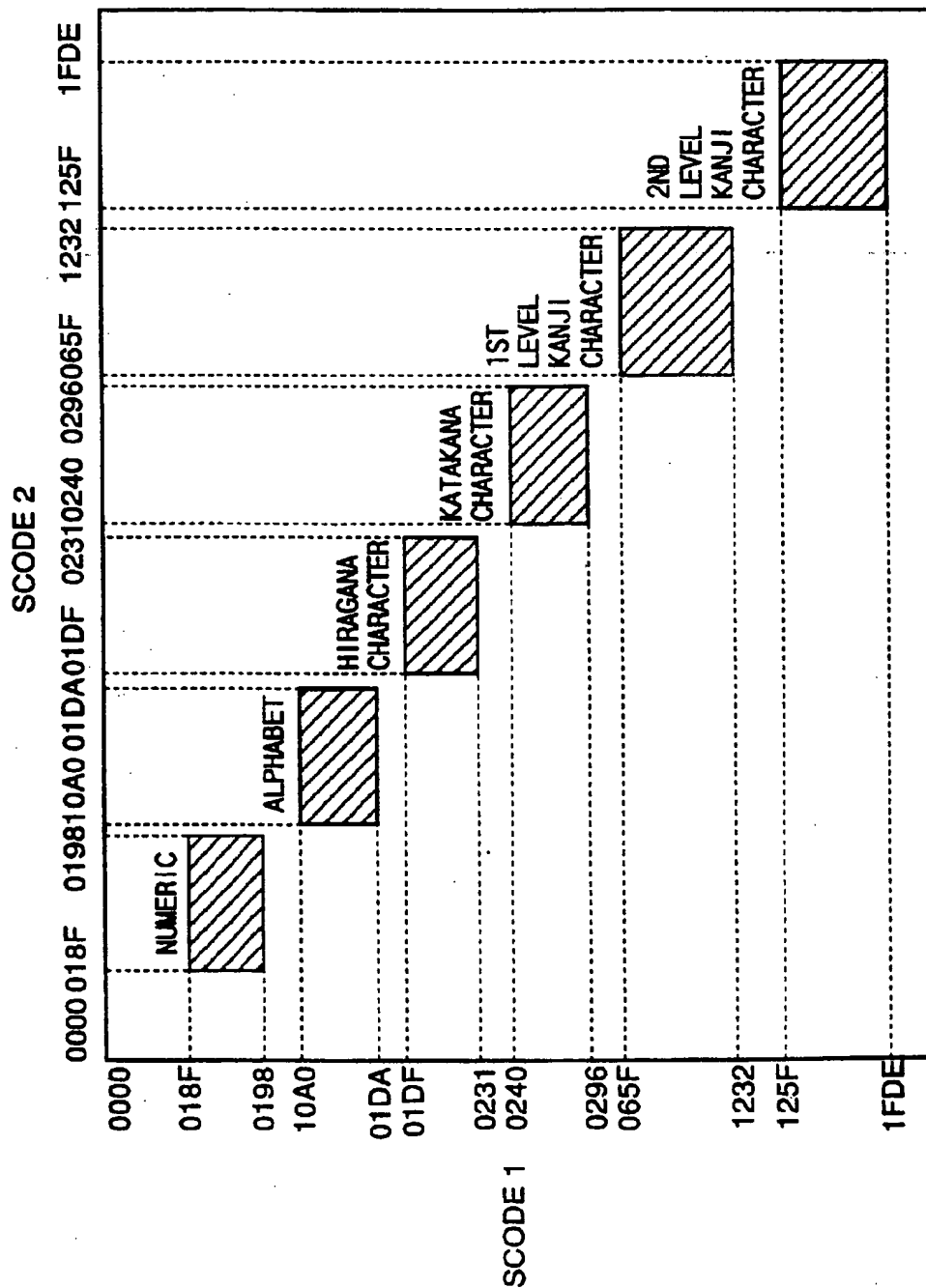


FIG. 57





## FIG. 58

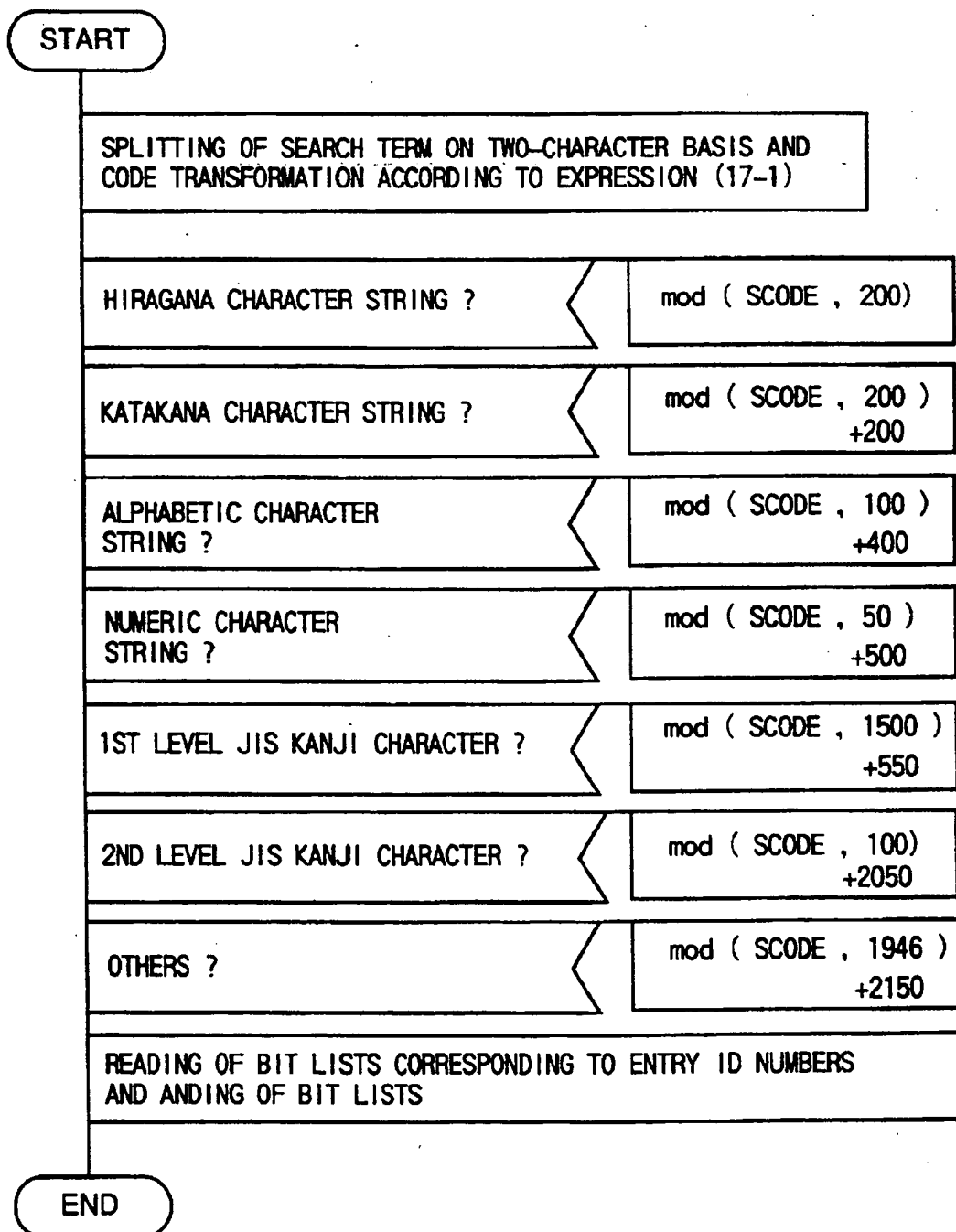
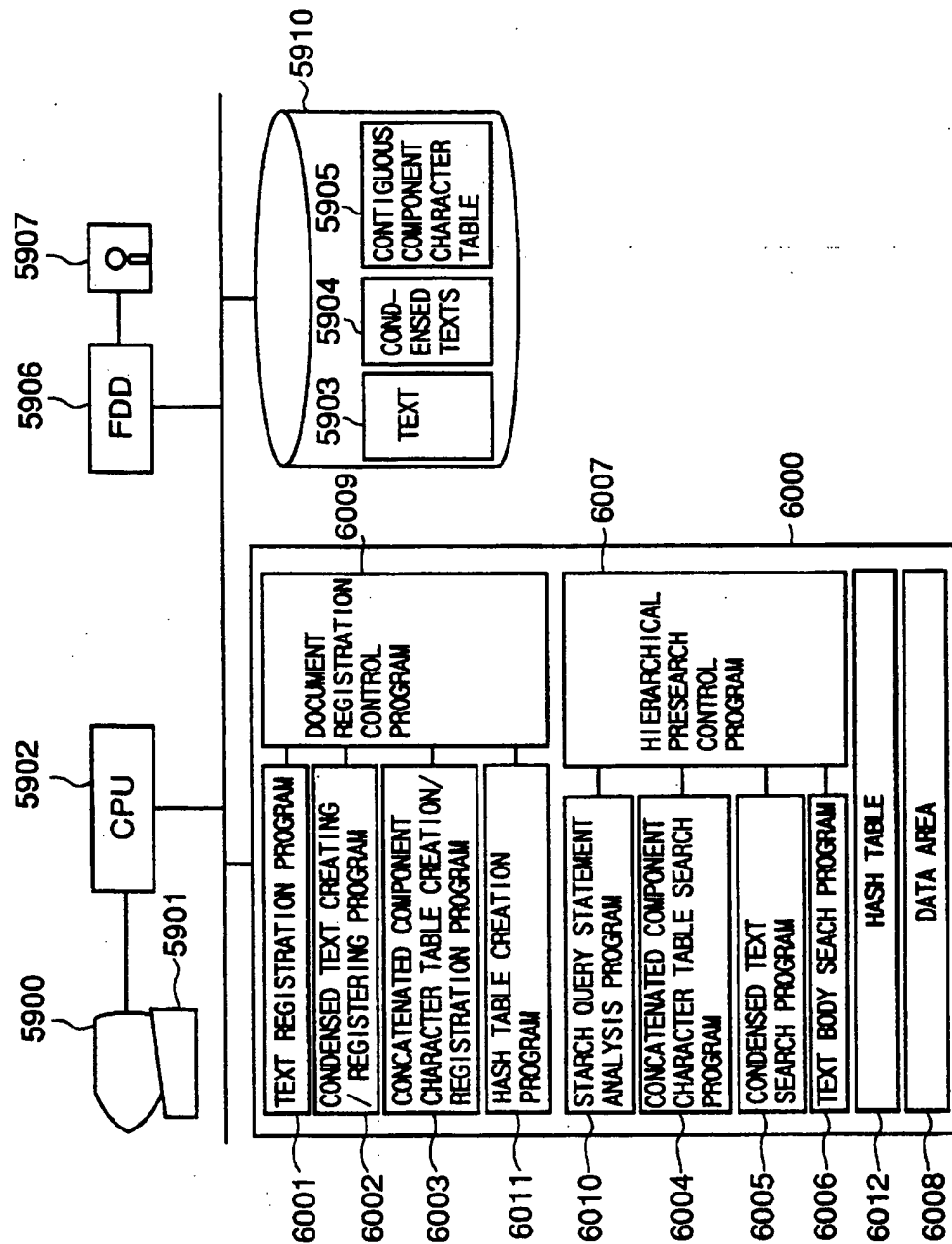
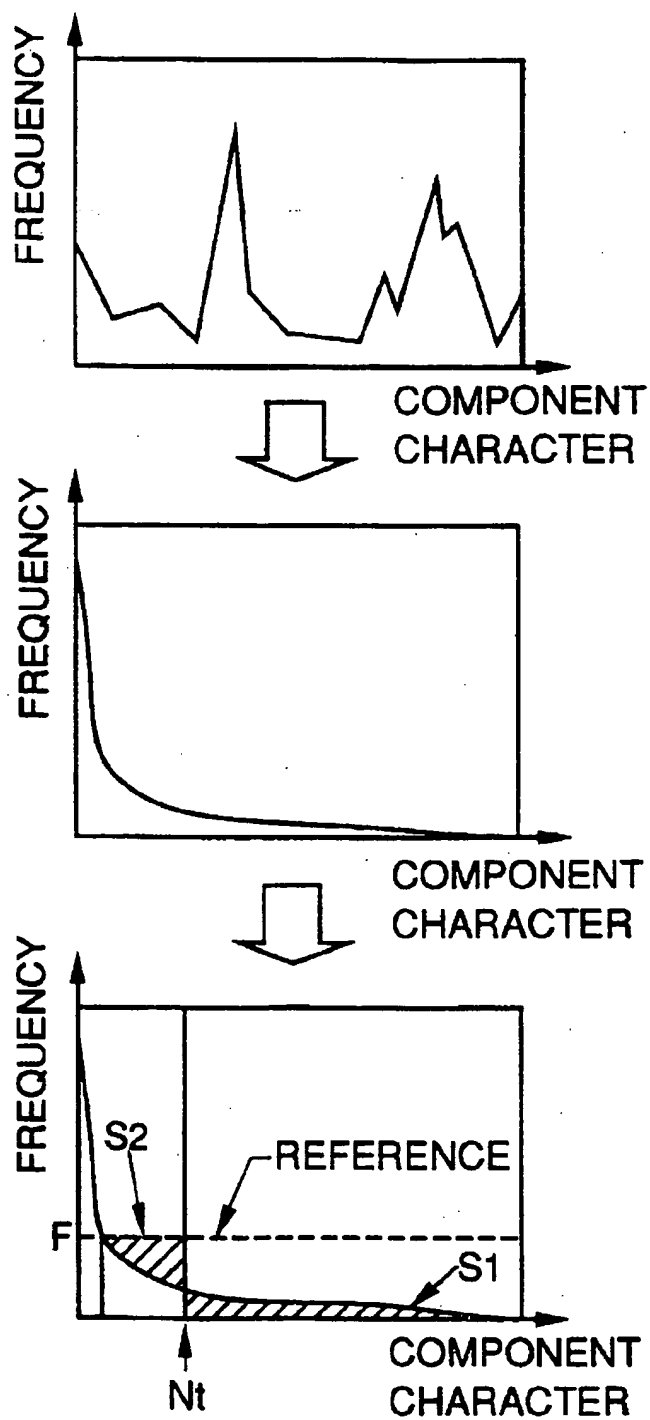


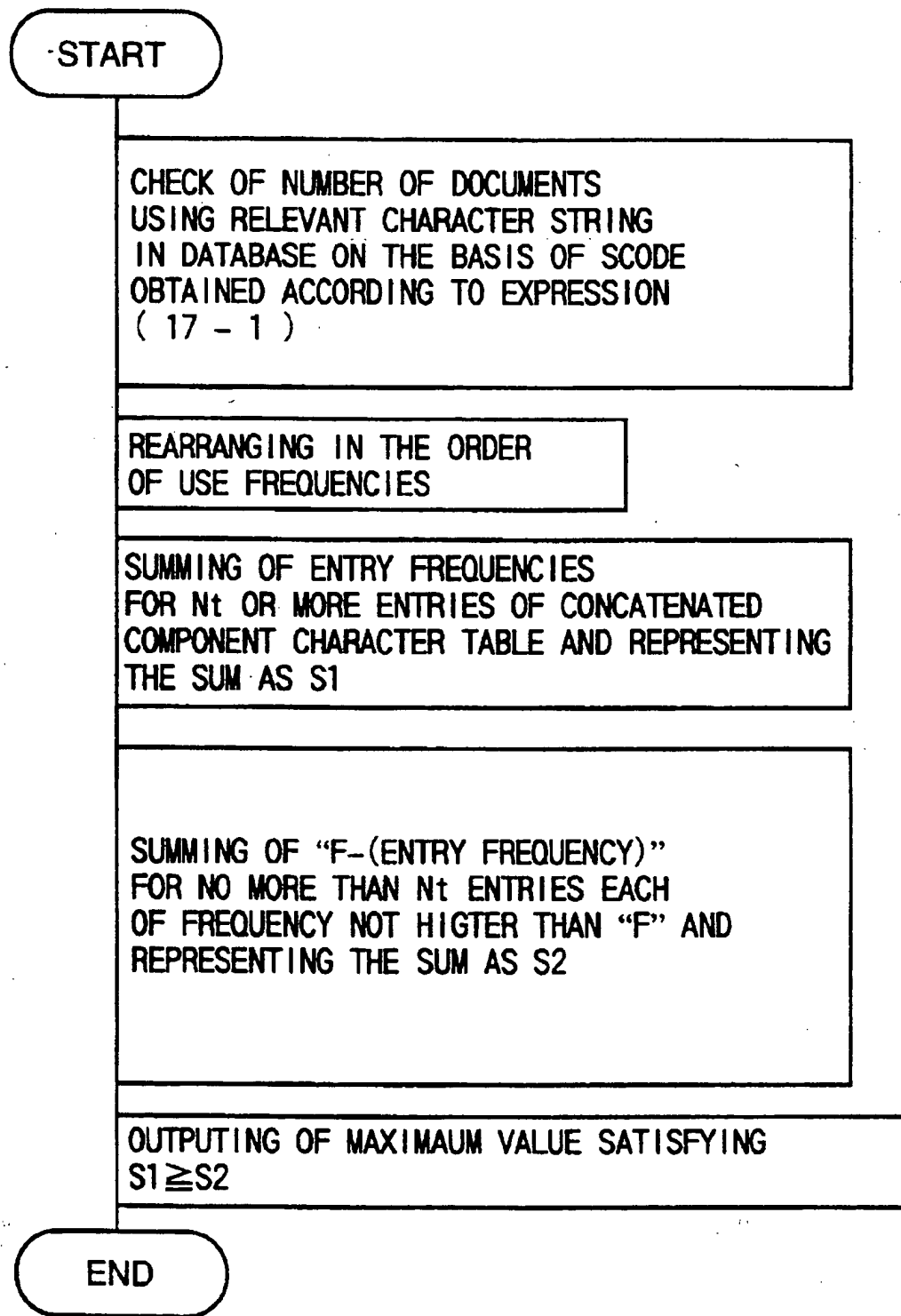
FIG. 59



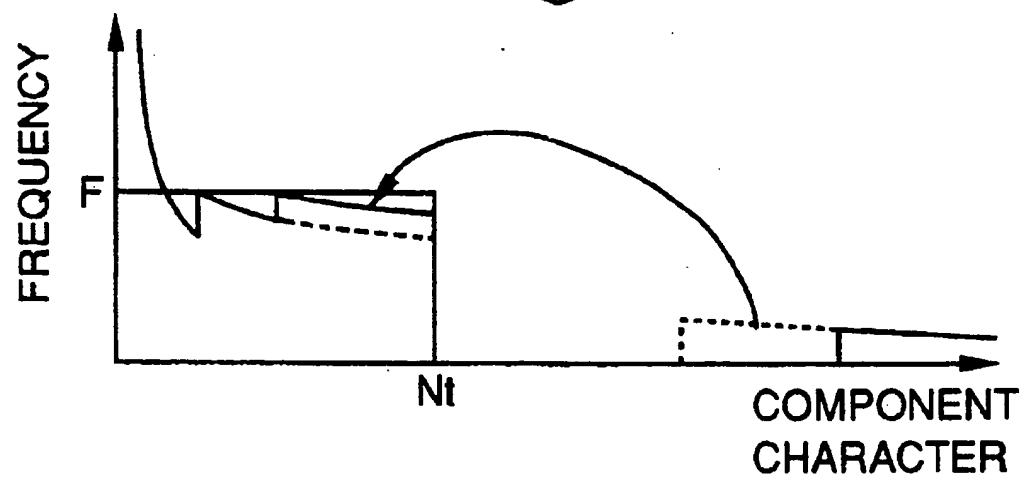
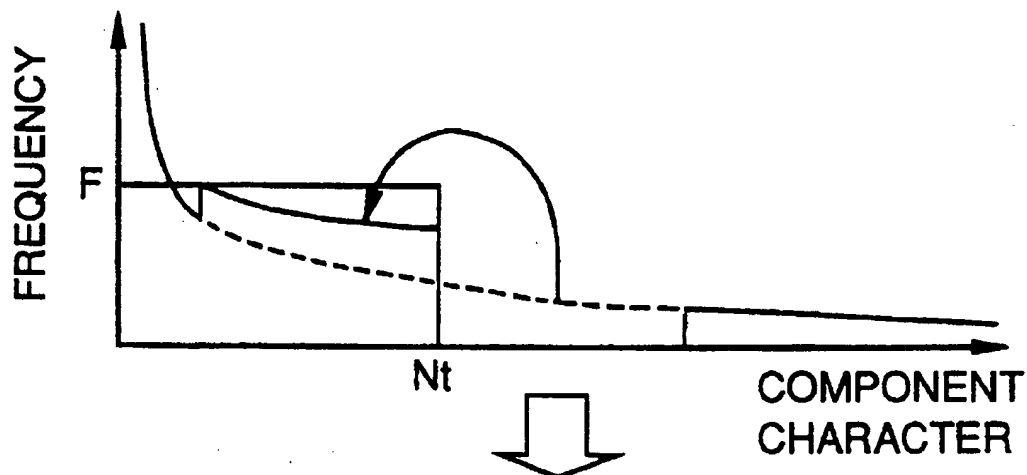
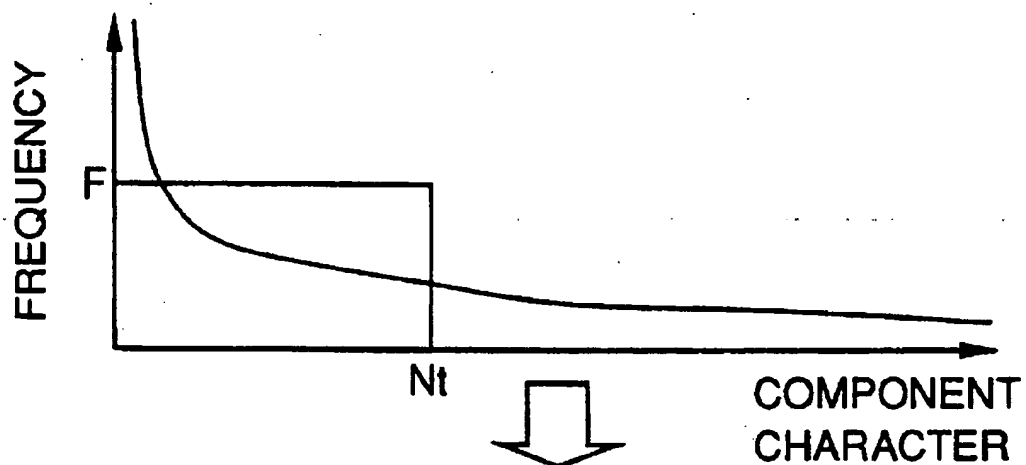
## FIG. 60



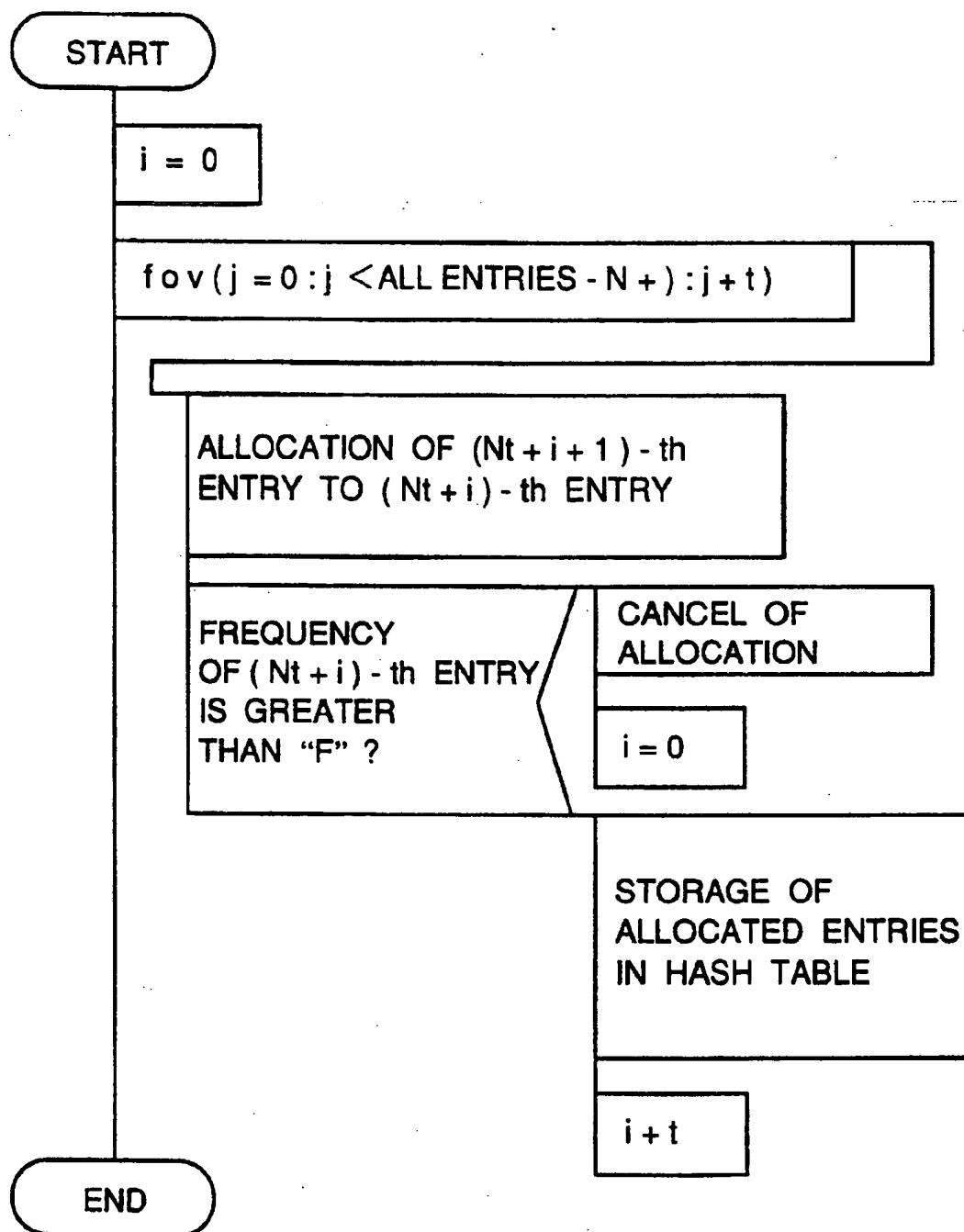
## FIG. 61



## FIG. 62



## FIG. 63



## FIG. 64

SCODE1	SCODE2	HASH ENTRY
.	.	.
(096F0A8C)	H	4032
(096F0A8D)	H	167
(096F0A8E)	H	1680
「経済」 → (096F0A8F)	H	34
(096F0A90)	H	2687
(096F0A91)	H	2948
(096F0A92)	H	862
.	.	.
.	.	.

## FIG. 65

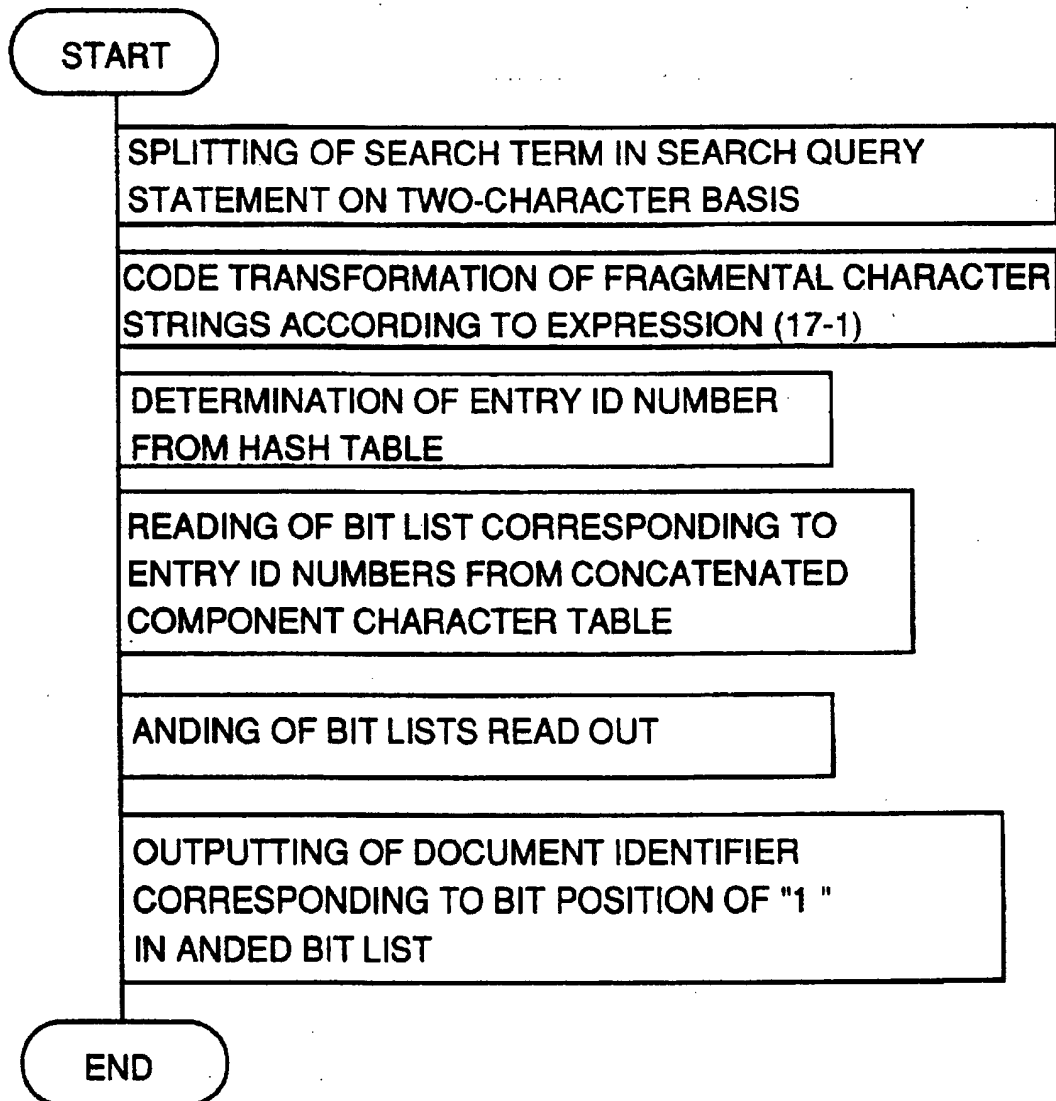
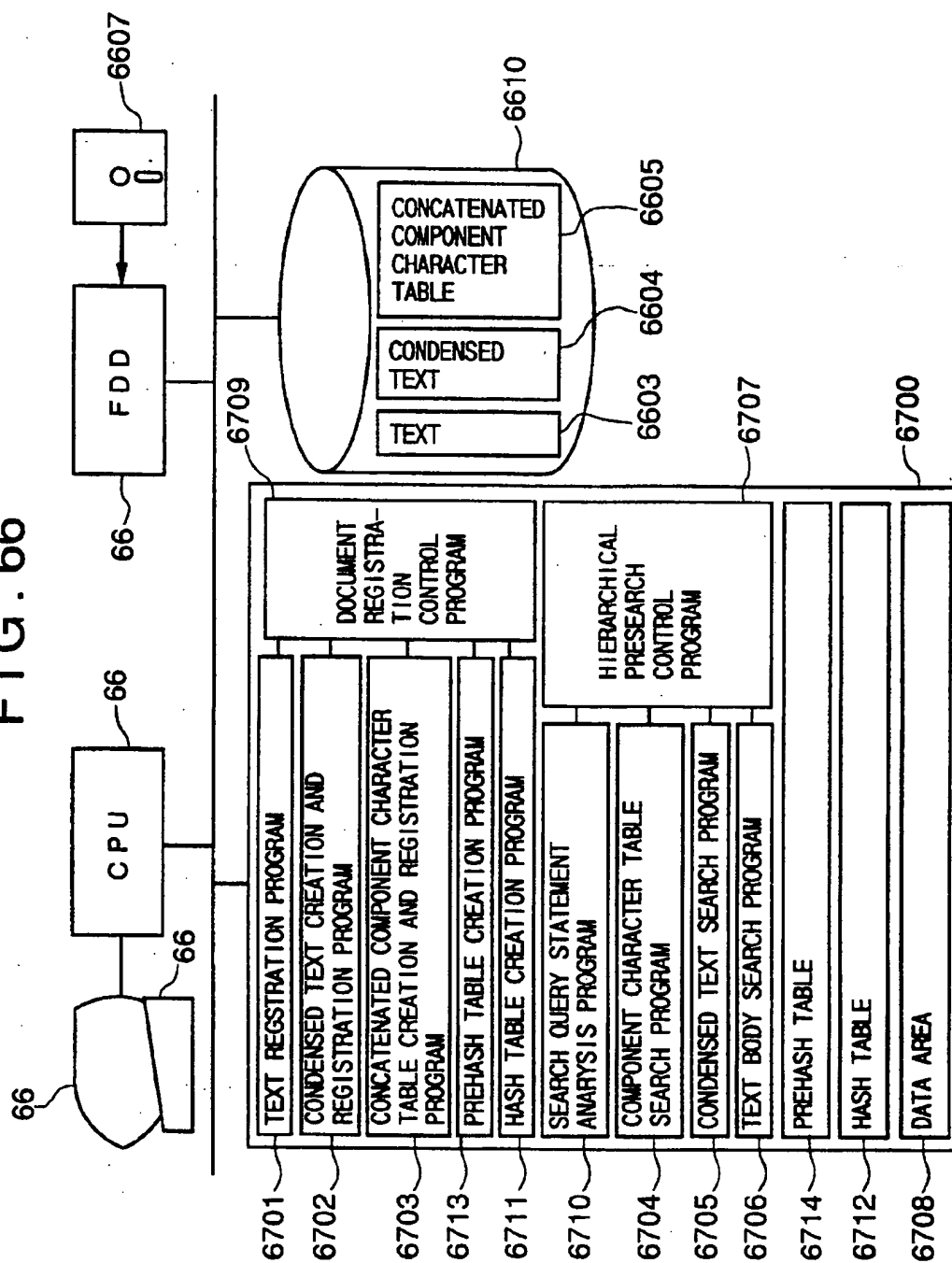
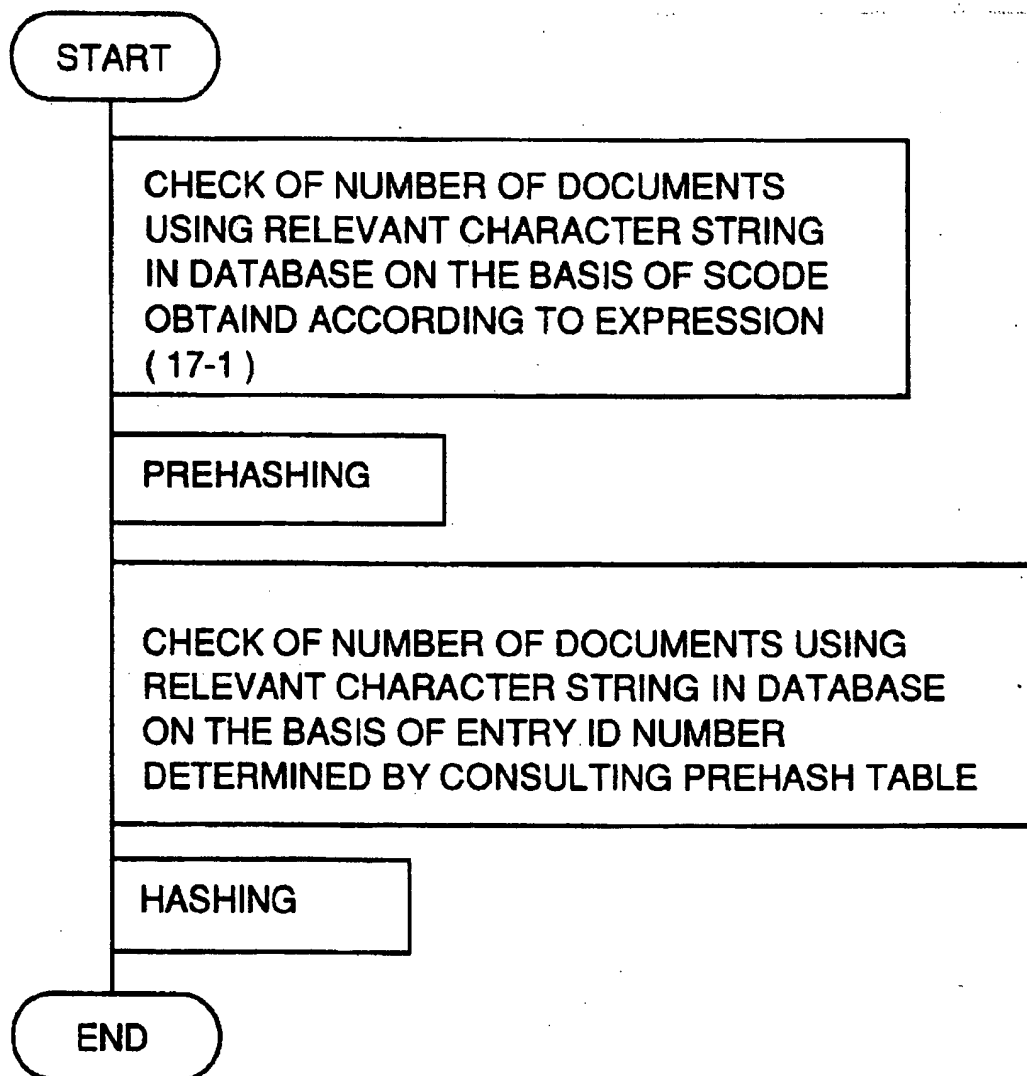




FIG. 66



## FIG. 67



# FIG. 68

